

An investigation of latency prediction for NoC-based communication architectures using machine learning techniques

SILVA, J, KREUTZ, M., PEREIRA, M. and DA COSTA ABREU, Marjory
<<http://orcid.org/0000-0001-7461-7570>>

Available from Sheffield Hallam University Research Archive (SHURA) at:
<https://shura.shu.ac.uk/25392/>

This document is the Accepted Version [AM]

Citation:

SILVA, J, KREUTZ, M., PEREIRA, M. and DA COSTA ABREU, Marjory (2019). An investigation of latency prediction for NoC-based communication architectures using machine learning techniques. Journal of Supercomputing, 1-19. [Article]

Copyright and re-use policy

See <http://shura.shu.ac.uk/information.html>

An investigation of latency prediction for NoC-based communication architectures using Machine Learning techniques

Jefferson Silva · Marcio Kreutz · Monica Pereira · Marjory Da Costa-Abreu

Received: date / Accepted: date

Abstract Due to the increasing number of cores in Systems on Chip (SoCs), bus architectures have suffered with limitations regarding performance. As applications demand higher bandwidth and lower latencies, buses have not been able to comply with such requirements due to longer wires and increased capacitance. Facing this scenario, Networks-on-Chip (NoCs) emerged as a way to overcome the limitations found in bus-based systems. Fully exploring all possible NoC characteristics settings is unfeasible due to the vast design space to cover. Therefore, some methods which aim to speed up the design process are needed. In this work, we propose the use of machine learning techniques to optimise NoC architecture components during the design phase. We have investigated the performance of several different ML techniques and selected the Random Forest one targeting audio/video applications. The results have shown an accuracy of up to 90% and 85% for prediction involving arbitration and routing protocols, respectively, and in terms of applications inference, audio/video achieved up to 99%. After this step, we have evaluated other classifiers for each application individually, aiming at finding the adequate one for each situation. The best class of classifiers found was the Tree-based one (Random Forest, Random Tree, and M5P) which is very encouraging and it points to different approaches from the current state of the art for NoCs latency prediction.

Jefferson Silva
Campus Universitário Lagoa Nova
Tel.: +55 84 988262814
E-mail: jefferson@ppgsc.ufrn.br

Marcio Kreutz
UFRN, E-mail: kreutz@dimap.ufrn.br

Monica Pereira
UFRN, E-mail: monicapereira@dimap.ufrn.br

Marjory da Costa-Abreu
UFRN, E-mail: marjory@dimap.ufrn.br

Keywords network-on-chip · machine learning · latency prediction · design space exploration.

1 Introduction

The evolution of lithography process has driven an increase in circuit integration leading to more complex architectures and higher development costs, allowing that multiple circuits can be integrated into a single package, which can be also known as System on Chip (SoC) [1, 2]. This level of integration has led multiple cores in an individual silicon wafer, creating Multiple Processors Systems on Chip (MPSoCs) [3]. As the number of cores increased in a system, we have seen an increased in the possible number of communications among them. In order to improve performance, energy usage and scalability issues, Networks on Chip were employed, which use some elements from Ethernet networks, such as routers and links [4] to implement communications inside SoCs.

Since each NoC internal router component can be implemented in many different ways, finding optimised architectures rely on the understanding of a vast design space. Architectural options for NoCs stand on several possible different implementations which can range from buffering, scheduling, routing and flow control, leading to an NP-hard combinatorial problem [5]. In such a scenario, optimal architectures could only be obtained by testing each one of all possible configurations, which is computationally unfeasible. However, since applications constraints could be known before hand, optimised architectures - the ones complying with application constraints - are usually satisfactory. Even so, in order to get there, many different possible behaviors would need to be tested. This task could be accomplished using meta-heuristics, such as Genetic Algorithms [6], by traversing the design space smartly enough to satisfy the constraints while avoiding testing all architectural scenarios. Still, each architecture found by the heuristic must be simulated to verify how far it is from the desired constraints. Usually, constraints are expressed in terms of performance, energy consumption, area.

Traditional approaches rely on simulations to run the architectures and get, for instance, performance figures. However, precise results (at cycle level) could only be achieved when simulations run at lower abstraction levels, which in turn, can take too much time, due to the high computation effort needed to simulate all the components at each clock cycle. In order to overcome this situation, many approaches rely on simulations performed at higher abstraction levels. Although this brings faster simulations times, results could be compromised due to lower accuracy.

Nowadays a trend is to avoid simulations using a variety of methods, such as Analytical Modeling [7] and Machine Learning [8–10]. This tendency is satisfactory because it allows an improvement in accuracy (avoiding not obeying restrictions of time, required area or energy consumption) and a speedup in design time. Nonetheless, in all cases, designers seek to increase accuracy and

to reduce the error rate. Despite this goal, NoCs have a non-linear behavior, which compromises the accuracy, and it is a challenge by itself. The accuracy loss could cause resource waste because the designer will use more resources than necessary aiming to obey the deadlines. Thus, it is a possible disadvantage of predictions.

Following this trend, this work proposes the usage of machine learning methods to predict latency values based on NoC architecture details. The choice by Machine Learning methods is useful for applications where the data is difficult to model analytically and NoC has this characteristic [11]. According to Ogras, Hu and Marculescu [12] and Qian et al. [13], there are at least three main challenges in NoC design:

1. Assumptions in current queuing models (analytical models) are very tight, not supporting a wide range of traffic patterns;
2. Efficient resources management strategies;
3. Scalability challenge in NoC simulations. The authors state that simulating NoCs with more than one hundred routers demand huge computation effort, which leads to high simulation times.

Regarding the first item above, this work supports seven NoC characteristics and two application attributes to reliably predict latency, enabling the verification of more NoC configurations. The attributes are: topology, size, routing protocol, virtual channel, input buffer depth, output buffer depth, arbiter, number of packets, and required bandwidth. These last two attributes are provided by application. It impacts directly in challenge 2, as being able to explore diverse architectural configurations allows the adoption of more efficient resource management strategies. Concerning the last challenge cited above, scalability in simulations, our predictor was trained with up to 144 routers. Still, it is possible to increment this number even further, because, with our approach, it is not necessary to simulate the NoC behavior on every single iteration of the classifier, only what is required to feed the classifier. The main goal of this work is applying the ML method, to analyse its ability to predict (correctly or not) and investigate the reason behind the accuracy predicted.

The paper is organised as follows: In Section 2, we present the related works. Section 3 explains the proposed framework, which is discussed in more details in Section 4. Section 5 presents the obtained results, followed by conclusions and future directions.

2 Related Works

There are two major options to avoid having to perform all possible simulations during design space exploration:

- the use of Analytical Models; and
- the use of Machine Learning techniques (ML).

This section presents some related works and their contributions to this problem.

2.1 Analytical Models

Many NoC latency models use the queuing theory, and it implies in restrictions, such as packet length and traffic follows a Poisson distribution [14].

Feng, Ge, and Wu presented an analytical model that assumes infinite buffers to predict latency in 3D NoCs [15]. It also supports multi-application mapping. It, initially, uses a Genetic Algorithm to realise the task mapping and, the resulting mapping is applied to an analytical model to verify if it obeys the designer restriction for each application. Task Graphs For Free (TGFF) [16] was used to generate synthetic graph applications, and the result of Genetic Algorithm (GA) was simulated using NoC simulator Nirgam. This model can save power consumption up by 21% and decrease the latency up to 17% when compared with random mapping.

Quian et al. developed a model that uses a G/G/1/K queuing model that supports heavy traffic, finite buffer, arbitrary packet length, and round-robin arbiter, under synthetic traffic patterns, the model achieved less than 12% of error in prediction the network saturation point. Using real applications, the proposed model can achieve 91.4-97.9% of accuracy, comparing the required time to simulate (using Booksim simulator) and the speedup was 70x over the simulation [17].

Bhattacharya and Jha created a model that uses a G/G/1/K and M/G/1/K that can reflect the influence of buffers size, number of virtual channels, and flit size. The model was evaluated using gem5 and GARNET simulator with PAR-SEC benchmark suite. The obtained results show that the number of virtual channels and flits contention are inversely proportional; besides it, the authors cannot demonstrate the influence of buffers size in packet latency. Comparing the simulations results with predictions, the accuracy was superior to 85% in all scenarios [18].

Kurihara and Li proposed a model that supports four topologies, mesh, torus, hypercube, and metacube. The goal is to estimate the ratio cost-performance for these topologies. In relations to other NoC characteristics (buffers size, kind of arbiter.), anyone was considered in the model. According to the obtained results, torus and metacube topologies only achieve better cost-performance when the number of cores overcomes 100 units [19].

Fischer, Fehske, and Fettweis also created a model based on queuing theory, differently to the previously cited papers, they did not take into account information such as topology or routing scheme and assumed two characteristics: infinite buffers in input channels and that the traffic follows the Poisson distribution. It considers an NoC such a hierarchical structure, in this way, is possible to split in three steps: local arrives rate, forwarding probabilities, and path delays. Thanks to this division, the model can calculate each router

individually and after sum all. The accuracy, in obtained results, achieved 96% [20].

Fresse et al. dimensioned an NoC from mathematical models, to compare, they synthesised in FPGA board and analysed the results. An explicit restriction was the use of one fixed configuration during the experiments (Mesh 2D, two virtual channels, round robin arbiter, and XY routing algorithm). It limited the analysis of obtained results, and the authors affirmed that changing the kind of arbiter, the accuracy falls. In the worst case, this solution reached up by 88% of accuracy. Another characteristic is the low accuracy for small NoCs, and it allows us to conclude that the model can not represent all situations adequately (do not capture the non-linear behavior) [21].

According to the Quian et al., analytical models cannot represent the real applications, because they assume hard restrictions, trying to relax these assumptions, they developed a new model that uses a Support Vector Machine to improve the model accuracy [22]. However, besides it, the SVM fed the analytical model. This combination achieved a result of 10% better than the previous work for the same authors - only using a mathematical model.

As mentioned earlier, the goal is to predict the value for the desired metric and an accurate inference, the number and the possibilities for each attribute are essential because it allows the model to create a realistic representation of NoC. These quoted analytical models need tight restrictions and are not capable of representing real applications accurately. In this way, the proposed framework operates with a high degree of freedom about the requirements. About the time distribution of arrived packets, our solution was trained with a variety of packet time distributions. Other difference regards on the possibility to choose among several NoC topologies sizes. The proposed solution was trained with two topologies ranging from 4 to 144 routers. Each author adopted a set of information to feed his model and to predict the desired metric, but anyone of the quoted paper informed the use of seven and two NoC characteristics and applications, respectively. Thus, our approach supports a higher number of applications because the proposed model supports more options for each attribute than the presented in the quoted papers, such as routing protocol which this paper handles six different protocols.

2.2 Machine learning-based solutions

Machine Learning (ML) techniques are commonly used for design space exploration purposes. For instance, Chen et al. use an ML technique to speed up networks up to 8 routers; they used RankBoost method to find the best configuration and creates a configuration rank to choose the best one, based on a defined metric threshold. However, this work did not focus on NoC DSE, but in CPU DSE [23].

Few works use ML methods to focus on regular NoC topologies, while the opposed happen for irregular ones. Specifically, about irregular topologies, Chou and Marculescu created a platform focused in user experience to explore

the design space using Machine Learning techniques to cluster the traces from similar users and, for each cluster, an algorithm for automatically architecture generation. This algorithm in a second phase implements an analytical model and takes into consideration many characteristics such as architecture template (resources of computation, communication, protocols), applications specification (task graph, deadlines, power consumption), and user traces (restrictions about latency, power consumption) [24].

In this direction, Zhang et al., used a congestion matrix to predict the worst packet latency in NoCs. According to the authors, the worst-case traffic model can adequately reflect the dynamic behaviors of packet switching networks. To find the best configuration (the solution supports seven attributes) for the performance they used a local search algorithm to explore NoC configurations; the platform also provides an area and power consumption estimations, generated using ORION [25] and DSENT [26] tools. The results showed that the bigger is the NoC size, lesser is the solution accuracy. They also presented accuracy results over injection rates. In this case, for lower injection rates, such as from 0.1 to 0.3, the accuracy found was less than 70% [27].

Aiming at obtaining an adequate NoC configuration, Ayari et al. developed a hypervolume-based approach to refining the DSE. They used as metric the load variation, power consumption, and communication costs. Because it is a multi-objective solution, the authors used a Reduced Pareto front (RPF) to facilitate the designer solution choice. The proposed approach had two stages: first, Pareto optimal solutions are clustered using K-means algorithm, and, as the second stage, a subset of solutions that maximises the hypervolume is selected through a genetic algorithm [28].

Different from the other approaches, our solution supports a significant number of routing and arbiters implementations. This is important since Fresse et al. argue that these characteristics may lead to improvements on the prediction errors due to the nonlinearity behavior of communications [21]. Comparing with irregular topologies, our work presents some advantages, such as minor design cost and, hence, higher possibility of reuse. Minor design cost happens because the architectures found can be shared among several applications. Contrasting with other approaches, our work is more flexible, not being tied to specific network sizes or target applications.

3 Proposed approach

This work aims to expand the use of classifier techniques to NoCs architectural characteristics in order to maximise its overall accuracy. The goal of this classifier is to predict average latency values based on NoCs and traffic pattern characteristics. For reaching this goal, we defined a sequence of work as presented in Figure 1.

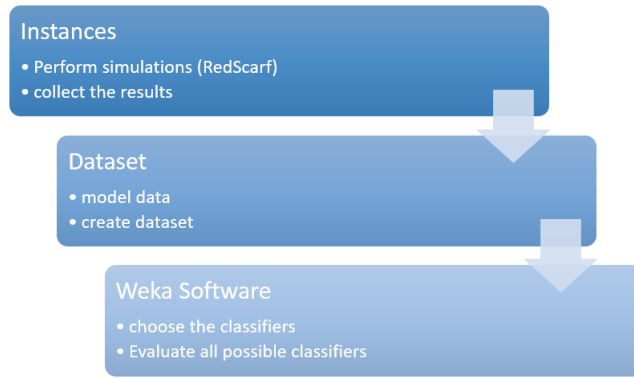


Fig. 1: Sequence of this approach.

Figure 1 defines the workflow. The first step is to perform simulations and to gather the results. In our case, we stopped this phase when we aggregated 496 instances. The data follows the requirements of Weka tool dataset. In the sequence, we identified all classifiers which operate in a regressive way implemented in the tool, and we evaluated all of them. We present more details for each step in the next subsections.

3.1 Machine Learning classifiers

Seeking out to predict accurately, the solution uses a regressive supervised method. We evaluated all regressive classifiers already implemented in Weka software. The list is presented below:

- **Isotonic Regression:** Selects the attribute that results in the lowest squared error;
- **Linear Regression:** It works by estimating coefficients for a line or hyperplane that best fits the training data;
- **MLP Regressor:** Trains a Multilayer Perceptron (MLP) with one hidden layer by minimising the given loss function plus a quadratic penalty with the Broyden-Fletcher-Goldfarb-Shanno method;
- **MLP:** An MLP consists of at least three layers of nodes. Except for the input nodes, all nodes using a nonlinear activation function. MLP uses a technique called backpropagation for training;
- **RBF Network and RBF Regressor:** An RBFN performs classification by measuring the inputs similarity to examples from the training set. When we want to classify a new input, each neuron computes the Euclidean distance between the input and its prototype;
- **SMOreg:** Implements the Support Vector Machine (SVM) for regression;
- **IBk:** K -nearest neighbours classifier. It uses the Euclidean Distance to calculate the distance between the instance and k neighbours;

- **M5Rules:** Generates a series of M5 trees, where only the "best" (highest coverage) leaf/rule is retained from each tree. At each stage, the instances covered by the best rule are removed from the training data before generating the next tree;
- **M5P:** This technique works as ordinary decision trees with linear regression models at the leaves that predict the value of observations that reach the leaf. The nodes of the tree represent variables and branches represent split values;
- **Random Forest:** The Random Forest algorithm consists of a random collection of decision trees. Hence, this algorithm is just an extension of the decision tree algorithm.

In order to use a well-known platform for others to be able to replicate our experiments, we have used the Weka¹ software, version 3.8. Each dataset represents an application. The statistical test Student t-test was applied to validate the results [29].

All needed simulations were performed on RedScarf Simulator [30]. It operates in Register Transfer Level (RTL) and has high accuracy, also has support to manipulate seven NoC attributes. Models for latency and flow rate follow the model's Dally and Towles [31].

3.2 Datasets

Audio/Video dataset contains 496 instances, each instance mapped to a specific NoC configuration plus two specific information about the application, precisely number of packets and required bandwidth, and the latency's value. Traffic generation was based in [32], and the communication distribution uses four models: Bit-reversal, Perfect Shuffle, Butterfly, and Transpose Matrix.

An instance has a network size between 2x2 to 12x12. These sizes were needed because the nonlinear variation in some circumstances, such as latency using nondeterministic routing algorithm, due to the occurred contention and possibility of deadlocks. Each instance represents a performed simulation on RedScarf simulator.

In addition to NoC attributes, two pieces of information about the application traffic pattern were used, the number of packets (NoP) and required bandwidth (RB). For NoP, values were 128, 1024, and 8192 per communication flow and for RB, 64, 512, and 1024 Mbps. All experiments were performed using all values for NoP and RB. This range is necessary to improve the accuracy of the classifiers.

There are three other datasets, one for each tested application: signaling, read/write, and block transfer. They contain 1,488 instances each. This number was reached varying NoC and application parameters. Similar to audio/video dataset, the same NoC values to attributes were used. Besides, applications characteristics were expanded. Number of Packets were evaluated

¹ <https://www.cs.waikato.ac.nz/ml/weka/>

with 128, 1024, 2048, and 8192; required Bandwidth was used four values: 64, 512, 1024, and 4096. The intention of using wide datasets was to achieve better training for classifiers.

All configurations for all characteristics gave us the total number of possibilities for a NoC, as showed in Table 1.

Characteristic	Number of possibilities
Topology	2
Size	15
Routing Protocol	6
Virtual Channels	33
Input Buffer Depth	28
Output Buffer Depth	33
Arbiter Type	4
Total	21,954,240

Table 1: All possibilities to NoC Design Space Exploration in this work.

By analysing the data on Table 1, one can conclude that there is a massive number of possibilities for each NoC and, hence, a high number of simulations to deplete all options. Assuming each simulation executing in ten seconds, 60,984 hours will be necessary to simulate everything (sequential simulations). Therefore, the DSE is huge and is almost unfeasible to evaluate all possibilities in an affordable time. Besides, there is another problem: each change made may impact in other attributes, leading designers to be forced to perform still more simulations and analysis to know the NoC final performance, which requires even more time. Based on this scenario, the values in Table 1 justify the use of AI methods for predictions.

Dataset was formatted according to Weka software ². Each experiment execution was repeated ten times, this value is a suggestion of Weka developers. Information about the RedScarf simulator is presented in the next Section.

3.3 RedScarf Simulator

RedScarf simulator implemented using System-C language with supports RTL simulation and multi-threading operation [33]. It supports real-time and non real-time applications. Figure 2 exhibits RedScarf configuration.

Figure 2 presents a topology with size equals to 8x8 and its traffic following Butterfly distribution. Figure 2a specifies the created traffic, where the designer can configure the traffic characteristic, some attributes are: distribution pattern of traffic, application (traffic class), number of packets, deadline

² <https://www.cs.waikato.ac.nz/ml/weka/downloading.html>

Traffic Edit

Source node: { **Multiple** }

Destination node(s): Butterfly

Specific addressing: Address: 0

Traffic class: RT0 - Signalling

Type of injection: Constant

Switching technique: Wormhole

Number of packets per flow: 4096

Deadline (ns): 0

Required bandwidth (Mbps): 1024

Message size (bits): 32

Idle time (ns): 0

Message interval (ns): 0

Function of probability: Normal

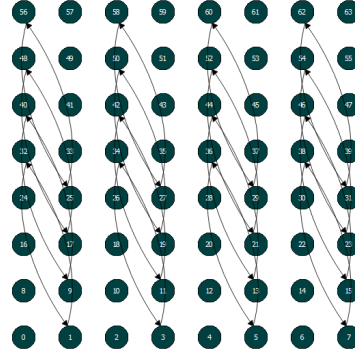
Std. deviation (% of req. bw): 1

Pareto - alpha on: 1.25

Pareto - alpha off: 1.90

Apply Cancel

(a) Traffic configuration in RedScarf.



(b) Mesh 8x8 topology with the specification of traffic among routers.

Fig. 2: Example of NoC configuration in RedScarf.

(for real-time applications), required bandwidth (also is known as packet injection rate), and message size. Figure 2b includes routers and its traffic (each arrow represents a communication flow).

The dataset was modeled aiming to represent the implemented attributes in RedScarf. The focus of simulator is performance simulation. Thus, it does not provide any information about the demanded area or power consumption. This tool offers a high precision because it operates in Register Transfer Level; despite this, it takes a long time to perform each simulation, which may become unfeasible a high number of simulations due the high demanded time.

Next Section will provide information about the Experiments.

4 Experiments and Results

For this research, the results were generated in two steps. First, it was necessary to simulate several NoC configurations, focusing on the average latency of packets. Each experiment was executed ten times. This value is the default in the simulator and is used by simulator's authors in their tests. These simulations were made using the RedScarf simulator running on a Debian 9 Linux.

After capturing the simulation results, the data was formatted in the Weka dataset file. Each attribute is related to a specific NoC characteristic. Thus, testing different NoCs account for changing any of the attributes (or even more than one) at a time. The next step was to evaluate all quoted classifiers listed in Section 3. Table 2 shows the student test result, and this statistical test analysed the accuracy information, not latency's result.

Table 2: t-test statistical test execution over all classifiers. Three classifiers obtained the same statistical significance.

Dataset	(1)	(2)	(3)	(4)	(5)
latencyPrediction	0.95 \pm 0.06	0.62 \pm 0.10 •	0.71 \pm 0.08 •	0.89 \pm 0.07 •	0.89 \pm 0.07 •

Dataset	(6)	(7)	(8)	(9)	(10)
latencyPrediction	0.51 \pm 0.24 •	0.90 \pm 0.06	0.89 \pm 0.06 •	0.85 \pm 0.09 •	0.93 \pm 0.05

◦, • statistically significant improvement or degradation

In Table 2, one can observe that three classifiers have obtained equivalent results. Table 3 presents the used configuration for each classifier; these values are the default in Weka software. Exploiting each hyperparameter of each technique is out of scope for this paper.

Table 3: Classifiers configuration used during the experiments.

- (1) functions.IsotonicRegression " 1679336022835454137
- (2) functions.LinearRegression '-S 0 -R 1.0E-8 -num-decimal-places 4'
- (3) functions.MLPRegressor '-N 2 -R 0.01 -O 1.0E-6 -P 1 -E 1 -S 1
-L functions.loss.SquaredError -A functions.activation.ApproximateSigmoid'
- (4) functions.MultilayerPerceptron '-L 0.3 -M 0.2 -N 500 -V 0 -S 0 -E 20
-H a' -5990607817048210779
- (5) functions.RBFNetwork '-B 2 -S 1 -R 1.0E-8 -M -1 -W 0.1'
- (6) functions.RBFRegressor '-N 2 -R 0.01 -L 1.0E-6 -C 2 -P 1 -E 1 -S 1'
- (7) functions.SMOreg '-C 1.0 -N 0 -I \"functions.supportVector.RegSMOImproved
-T 0.001 -V -P 1.0E-12 -L 0.001 -W 1\" -K \"functions.supportVector.Puk
-O 1.0 -S 1.0 -C 250007\"' -7149606251113102827
- (8) lazy.IBk '-K 1 -W 0 -A \"weka.core.neighboursearch.LinearNNSearch -A
\\\\\"weka.core.EuclideanDistance -R first-last\\\\\"' -3080186098777067172
- (9) trees.M5P '-M 4.0' -6118439039768244417
- (10) trees.RandomTree '-L 0 -M 1.0 -V 0.001 -S 1'

As can be seen in Table 2, three classifiers achieved the same performance. Thus, aiming to choose the best classifier, we analysed the Mean Absolute Error (MEA) as the second metric. Table 4 presents the results for this metric for the best classifiers.

Table 4 shows percentage of error for each of three classifiers. Random Tree method reached the minor error, without penalising accuracy.

Figure 3 presents an analysis of accuracy for each NoC size.

Figure 3 shows the accuracy of the Random Tree, Isotonic Regression, and SMOreg (Support Vector Machine for regressive problems) for several network sizes. Besides the high variation, from zero to near 1, the error rate

Table 4: Obtained values of Mean Absolute Error metric for each of three classifiers.

Classifier	Mean Absolute Error (%)
Isotonic Regression	21
SMOreg	13
Random Tree	9

Accuracy of the best three classifiers over network size

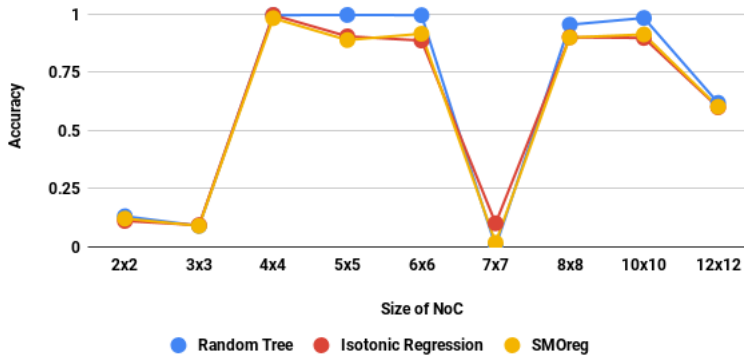


Fig. 3: Comparison between achieved accuracy using Isotonic Regression, Random Tree, and SMOreg for a range of NoC sizes.

sustained lower values. Analysing the accuracy for size seven, the classifiers almost missed all instances, but with a smaller error (less than 10%). In five sizes, out of 9 in total, the classifier overcame 95% of accuracy. This kind of behavior is expected because of the intrinsic behavior of NoCs, such as contentions. Comparing them - all of them using the hyperparameter presented in Table 3 - is possible to notice that the techniques had difficulty to infer the latency in some sizes. This allows us to conclude that the problem is not the adopted techniques by themselves rather the intrinsic behavior of NoCs, which is non-linear and may get contentions depending on application or traffic pattern distribution. For now on, all experiments were conducted using Random Tree classifier. Figure 4 shows the classifier accuracy for four applications.

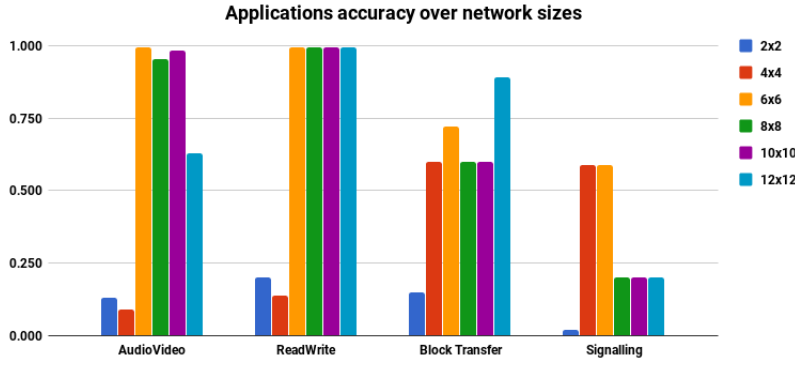


Fig. 4: Achieved accuracy for a range of NoC sizes using four applications.

As can be seen in Figure 4, read/write application achieved an accuracy higher than 99% from 6x6 to 12x12 networks. On the other hand, signaling application had the worst accuracy; only two NoC sizes achieved over 50% of accuracy. Audio/video application had similar results as shown in Figure 3.

For all applications, in 2x2 network size, the classifier reached up to 20% of accuracy. Despite the training process has the same number of instances for each network sizes, this result indicates that the Decision Tree used could not predict the latency values correctly for small NoCs. This could happen due to the nonlinear behavior of networks, imposing constraints for the learning process. Comparing the application distribution as presented in Figure 5 we can see some differences.

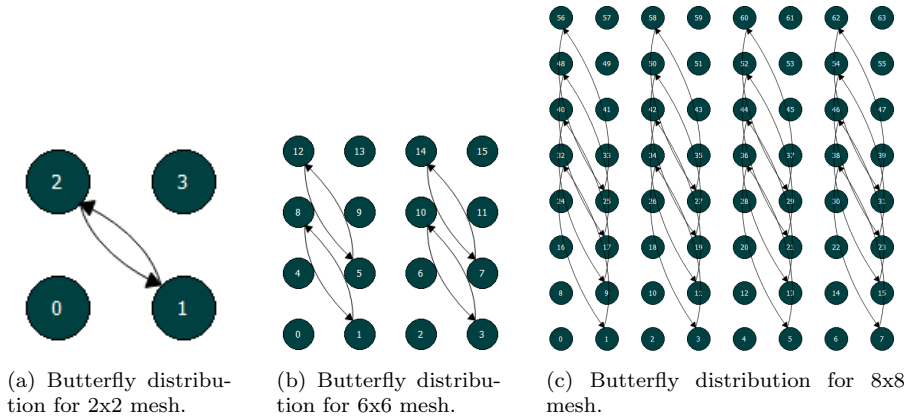


Fig. 5: Example of NoC applications using butterfly distribution, but with three different sizes.

Figure 5 shows three traffic patterns, one for each NoC size, 2x2, 6x6, and 8x8. In all cases, the application was mapped using butterfly distribution. Thus, the number of communications changes from two to thirty-two in Figure 5c. This change could impact directly in performance because it increases the possibility to occur contention, saturate the channels, for instance. Therefore, this NoC characteristic becomes more complex the task of the ML method to understand and generate a suitable model to a wide range of NoC sizes. Figure 6 presents the classifier accuracy for buffer sizes.

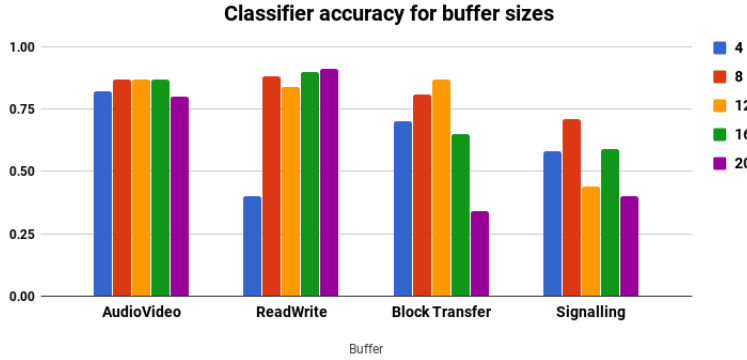


Fig. 6: Classifier accuracy for buffer sizes using four applications.

Figure 6 shows that the used classifier achieved, in two of five cases for signaling application, accuracy up to 45%. For block transfer application, buffer size equals to 16 and 20 had the worst accuracy, less than 70%, and 40%. In other cases, the accuracy overcame 70%. Read/write application achieved an accuracy greater than 80% in four scenarios. These results corroborate with the expected performance when the designer increases the buffer sizes: at a point, the performance increments proportionally, but, after reaching the limit, increase the buffer take to waste of resources. This behavior is unique for each application, each one has its resource usage. Therefore, the variation presented in Figure 6 is habitual. Taking Figure 2b as example and assuming the use of XY routing protocol [34], the router 25 forwards four communications: from router 32 to 1; from 40 to 9; from 48 to 17; from 56 to 25. Thus, this router needs the buffer to store the flits while it cannot forward ahead temporarily. However, on the other hand, router 0 does not need any buffer, because, in butterfly distribution, it does not receive any communication.

Based on previous results, it was analysed the possibility of usage of other classifier and Table 5 shows the best classifier for each application analysed.

Although each application requires a specific technique, all belongs to Decision Trees class. So, for the suggested approach, this class of methods is

Table 5: Best classifier for each application

Application	Method	CC	MAE	RAE
Signaling	Random Forest	0.9214	970.2602	14.3569%
Read/write	Random Tree	0.8753	1012.8454	14.0498%
Block transfer	M5P	0.906	466.6563	11.7552%

adequate. Figure 7 shows the accuracy using the best classifier for each application.

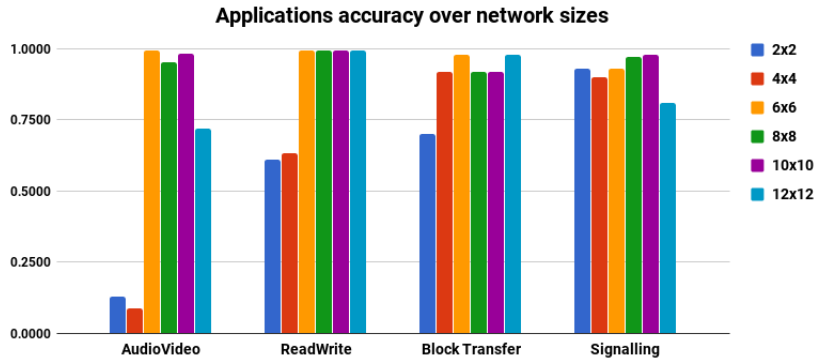


Fig. 7: Best classifiers accuracy to some NoC sizes using four applications.

Analysing each application separately, it is possible to notice that all three applications had an increase in accuracy. Signalling application overcame 70% of accuracy on overall. Worst scenario for Read/Write application was 2x2, but with the most adequate classifier, it overcame 55% of accuracy. Block Transfer application in five, out of 6 in total, overcame 92% of accuracy. Despite this improvement, Audio/Video application still obtained a lower result for sizes equal to 2x2 and 4x4. Thus, even the best classifier for this application, could not achieve high accuracy for both scenarios. Both sizes gathered the lowest accuracy for also Read/Write application, which demonstrates the lower requirement for communication and the disability of classifier to handle both sizes for these applications. However, for all applications, the accuracy increased for sizes equal or major than 6x6; it represents the classifier could modeling this behavior correctly.

Below, there is the Figure 8. It presents the accuracy for buffer sizes using the adequate classifier.

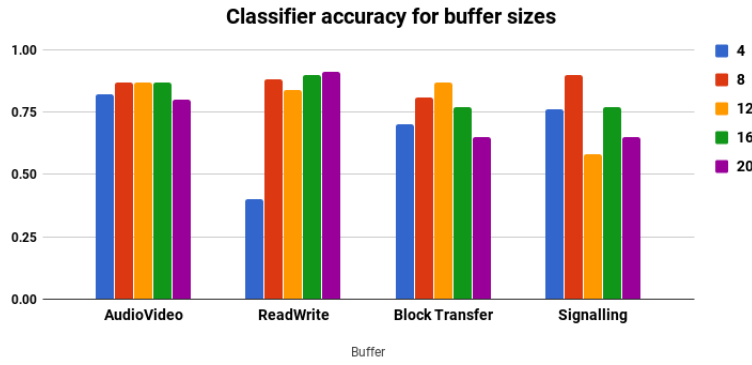


Fig. 8: Best classifiers accuracy for buffer sizes using four applications.

As can be seen, there was an increase in accuracy for three applications (all except Audio/Video) when adopted the best classifier for each application. For the signaling application, the results overcame 80% of accuracy at the worst scenario, and, for other scenario overcame 93% of accuracy. Block Transfer application had a similar growth, having like the worst case when buffer size is equal to 4 (less than 76% of accuracy). Read/write application had an accuracy higher than 50% in the worst case (also when buffer size is 4) and overcame 91% in four of five cases. Although the error rate appears to be significant, the error was always between 11% and 14%. In this way, it is possible to tune the values based on this error rate. Figure 9 presents the accuracy for all routing protocols supported.

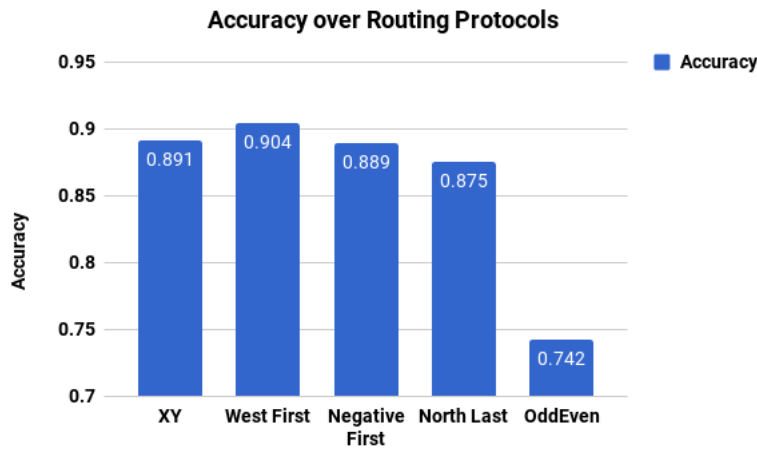


Fig. 9: Accuracy of classifier for each routing protocol supported.

By observing Figure 9, it is possible to notice that accuracy for OddEven protocol is less than 80%. This can be explained by the adaptive nature of the protocol that could lead to deadlocks, therefore making the prediction harder than for the other protocols. For the other four protocols, accuracy overcame 85%, which can be justified because the deterministic behavior of these routing protocols. Other similar situation is shown in Figure 10.

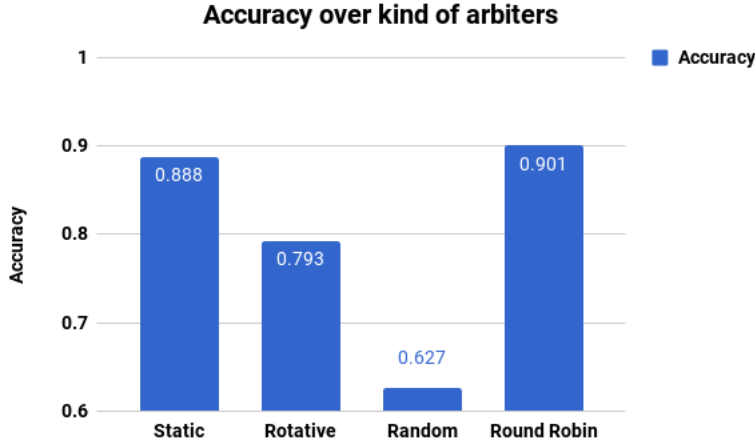


Fig. 10: Accuracy of classifier for each arbiter supported

Figure 10 walks in the same way that the previous image: classifier had difficulty in predicting adaptive or random situations. More specifically, random arbiter was the worst scenario, reaching less than 70% of accuracy. However, for the other types of arbiters, the accuracy overcame 79%.

An option to improve the accuracy for all cases is the use of a committee of classifiers. In this direction, Committee classification units could combine the outputs of two or more classifiers to generate a unique result [35].

5 Conclusions and Future Works

This work presented research about latency prediction for Networks on Chip communication architectures, based on the characteristics of the architectures and applications. Artificial Intelligence techniques were employed to build solutions able to predict latencies figures up to 99% in accuracy.

Networks on Chip were evaluated with mesh topologies ranging between 4 and 144 nodes. The experiments were conducted in two steps, taking simulation results to train the predictor.

Random Tree was used as Latency Predictor for nine attributes. The instances were collected from the RedScarf simulation tool. In total, four real applications were analysed.

Since good accuracy could be achieved by the predictor, it can be tuned to help speeding up the design space exploration tools for NoC based communication architectures. This is an essential contribution of the paper due to the novel approach proposed to predict the average latency for distributed applications.

As future works, we intend to extrapolate the experiments for other performance metrics and to analyse the usage of a committee of classifiers to improve even more the accuracy. Another possible direction relies on the use of heuristics to select the adequate characteristics for target architectures, aiming at reducing the needed resources and time to get optimised solutions.

6 Acknowledgments

This research was supported by High Performance Computing Center at UFRN (NPAD/UFRN).

References

1. C. A. Zeferino, Redes-em-chip: arquiteturas e modelos para avaliação de área e desempenho.
2. K. Jain, S. K. Singh, A. Majumder, A. J. Mondai, Problems encountered in various arbitration techniques used in noc router: A survey, in: Electronic Design, Computer Networks & Automated Verification (EDCAV), 2015 International Conference on, IEEE, 2015, pp. 62–67.
3. W. Wolf, A. A. Jerraya, G. Martin, Multiprocessor system-on-chip (mpsoc) technology, IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems 27 (10) (2008) 1701–1713.
4. T. Bjerregaard, S. Mahadevan, A survey of research and practices of network-on-chip, ACM Computing Surveys (CSUR) 38 (1) (2006) 1.
5. S. Kahruman-Anderoglu, A. Buchanan, S. Butenko, O. A. Prokopyev, On provably best construction heuristics for hard combinatorial optimization problems, Networks 67 (3) (2016) 238–245.
6. E. Elbeltagi, T. Hegazy, D. Grierson, Comparison among five evolutionary-based optimization algorithms, Advanced engineering informatics 19 (1) (2005) 43–53.
7. L. Ost, M. Mandelli, G. M. Almeida, L. Moller, L. S. Indrusiak, G. Sassatelli, P. Benoit, M. Glesner, M. Robert, F. Moraes, Power-aware dynamic mapping heuristics for noc-based mpsocs using a unified model-based approach, ACM Transactions on Embedded Computing Systems (TECS) 12 (3) (2013) 75.
8. Y. Z. Tei, M. N. Marsono, N. Shaikh-Husin, Y. W. Hau, Network partitioning and ga heuristic crossover for noc application mapping, in: Circuits and Systems (ISCAS), 2013 IEEE International Symposium on, IEEE, 2013, pp. 1228–1231.
9. H. R. Faragardi, R. Shojaei, N. Yazdani, Reliability-aware task allocation in distributed computing systems using hybrid simulated annealing and tabu search, in: High Performance Computing and Communication & 2012 IEEE 9th International Conference on Embedded Software and Systems (HPCC-ICES), 2012 IEEE 14th International Conference on, IEEE, 2012, pp. 1088–1095.

10. C. Wu, C. Deng, L. Liu, J. Han, J. Chen, S. Yin, S. Wei, A multi-objective model oriented mapping approach for noc-based computing systems, *IEEE Transactions on Parallel and Distributed Systems* 28 (3) (2017) 662–676.
11. V. Sze, Y.-H. Chen, J. Emer, A. Suleiman, Z. Zhang, Hardware for machine learning: Challenges and opportunities, in: 2017 IEEE Custom Integrated Circuits Conference (CICC), IEEE, 2017, pp. 1–8.
12. R. Marculescu, J. Hu, U. Y. Ogras, Key research problems in noc design: a holistic perspective, in: *Hardware/Software Codesign and System Synthesis, 2005. CODES+ISSS'05. Third IEEE/ACM/IFIP International Conference on*, IEEE, 2005, pp. 69–74.
13. Z. Qian, P. Bogdan, C.-Y. Tsui, R. Marculescu, Performance evaluation of noc-based multicore systems: From traffic analysis to noc latency modeling, *ACM Transactions on Design Automation of Electronic Systems (TODAES)* 21 (3) (2016) 52.
14. U. Y. Ogras, P. Bogdan, R. Marculescu, An analytical approach for network-on-chip performance analysis, *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 29 (12) (2010) 2001–2013.
15. G. Feng, F. Ge, N. Wu, Balancing 3d network-on-chip latency in multi-application mapping based on m/g/1 delay model, in: *Proceedings of the World Congress on Engineering and Computer Science*, Vol. 1, 2015.
16. R. P. Dick, D. L. Rhodes, W. Wolf, Tgff: task graphs for free, in: *Proceedings of the Sixth International Workshop on Hardware/Software Codesign.(CODES/CASHE'98)*, IEEE, 1998, pp. 97–101.
17. Z. Qian, D.-C. Juan, P. Bogdan, C.-Y. Tsui, D. Marculescu, R. Marculescu, A comprehensive and accurate latency model for network-on-chip performance analysis, in: *Design Automation Conference (ASP-DAC), 2014 19th Asia and South Pacific*, IEEE, 2014, pp. 323–328.
18. D. Bhattacharya, N. K. Jha, Analytical modeling of the smart noc, *IEEE Transactions on Multi-Scale Computing Systems* 3 (4) (2017) 242–254.
19. T. Kurihara, Y. Li, A cost and performance analytical model for large-scale on-chip interconnection networks, in: *Computing and Networking (CANDAR), 2016 Fourth International Symposium on*, IEEE, 2016, pp. 447–450.
20. E. Fischer, A. Fehske, G. P. Fettweis, et al., A flexible analytic model for the design space exploration of many-core network-on-chips based on queueing theory, in: *Proceedings of the Fourth International Conference on Advances in System Simulation*, ser. SIMUL, Citeseer, 2012.
21. V. Fresse, C. Combes, M. Payet, F. Rousseau, Noc dimensioning from mathematical models, *International Journal of Computing and Digital Systems (IJCDS)* 5 (2) (2016) 135–145.
22. Z. Qian, D.-C. Juan, P. Bogdan, C.-Y. Tsui, D. Marculescu, R. Marculescu, Svr-noc: A performance analysis tool for network-on-chips using learning-based support vector regression model, in: *Proceedings of the Conference on Design, Automation and Test in Europe*, EDA Consortium, 2013, pp. 354–357.
23. T. Chen, Q. Guo, K. Tang, O. Temam, Z. Xu, Z.-H. Zhou, Y. Chen, Archranker: A ranking approach to design space exploration, in: *Computer Architecture (ISCA), 2014 ACM/IEEE 41st International Symposium on*, IEEE, 2014, pp. 85–96.
24. C.-L. Chou, R. Marculescu, User-centric design space exploration for heterogeneous network-on-chip platforms, in: *Proceedings of the Conference on Design, Automation and Test in Europe*, European Design and Automation Association, 2009, pp. 15–20.
25. A. B. Kahng, B. Li, L.-S. Peh, K. Samadi, Orion 2.0: A fast and accurate noc power and area model for early-stage design space exploration, in: *Design, Automation & Test in Europe Conference & Exhibition, 2009. DATE'09.*, IEEE, 2009, pp. 423–428.
26. C. Sun, C.-H. O. Chen, G. Kurian, L. Wei, J. Miller, A. Agarwal, L.-S. Peh, V. Stojanovic, Dsent-a tool connecting emerging photonics with electronics for opto-electronic networks-on-chip modeling, in: *Networks on Chip (NoCS), 2012 Sixth IEEE/ACM International Symposium on*, IEEE, 2012, pp. 201–210.
27. M. Zhang, Y. Shi, F. Zhang, Z. Liu, Comrance: A rapid method for network-on-chip design space exploration, in: *Green and Sustainable Computing Conference (IGSC01)*, 2016 Seventh International, IEEE, 2016, pp. 1–8.

28. R. Ayari, M. Nikdast, I. Hafnaoui, G. Beltrame, G. Nicolescu, Hypap: A hypervolume-based approach for refining the design of embedded systems, *IEEE Embedded Systems Letters* 9 (3) (2017) 57–60.
29. P. Hsu, Contribution to the theory of student's t-test as applied to the problem of two samples., *Statistical Research Memoirs*.
30. E. A. SILVA, Redscarf: ambiente para avaliação de desempenho de rede-em-chip, Trabalho Técnico Científico de Conclusão de Curso, Universidade do Vale do Itajaí.
31. W. J. Dally, B. P. Towles, Principles and practices of interconnection networks, Elsevier, 2004.
32. L. Tedesco, A. Mello, D. Garibotti, N. Calazans, F. Moraes, Traffic generation and performance evaluation for mesh-based nocs, in: Proceedings of the 18th annual symposium on Integrated circuits and system design, ACM, 2005, pp. 184–189.
33. E. A. da Silva, D. Menegasso, S. Vargas, C. A. Zeferino, Redscarf: A user-friendly multiplatform network-on-chip simulator, in: 2017 VII Brazilian Symposium on Computing Systems Engineering (SBESC), IEEE, 2017, pp. 71–78.
34. M. V. Rao, T. R. Krishna, S. Siddhartha, K. S. Prathyusha, A load balance aware xy routing methodology for noc architectures, *Advances and Applications in Mathematical Sciences* 17 (1) (2017) 251–270.
35. M. Aksela, Comparison of classifier selection methods for improving committee performance, *Multiple Classifier Systems* (2003) 159–159.