# Optimizing Long Short-Term Memory Recurrent Neural Networks Using Ant Colony Optimization to Predict Turbine Engine Vibration

AbdElRahman ElSaid[a,*], Fatima El Jamiy[a], James Higgins[b], Brandon Wild[b], Travis Desell[a]

*University of North Dakota, Grand Forks, North Dakota 58202*

[a]*Department of Computer Science*
[b]*Department of Aviation*

**Abstract**

This article expands on research that has been done to develop a recurrent neural network (RNN) capable of predicting aircraft engine vibrations using long short-term memory (LSTM) neurons. LSTM RNNs can provide a more generalizable and robust method for prediction over analytical calculations of engine vibration, as analytical calculations must be solved iteratively based on specific empirical engine parameters, making this approach ungeneralizable across multiple engines. In initial work, multiple LSTM RNN architectures were proposed, evaluated and compared. This research improves the performance of the most effective LSTM network design proposed in the previous work by using a promising neuroevolution method based on ant colony optimization (ACO) to develop and enhance the LSTM cell structure of the network. A parallelized version of the ACO neuroevolution algorithm has been developed and the evolved LSTM RNNs were compared to the previously used fixed topology. The evolved networks were trained on a large database of flight data records obtained from an airline containing flights that suffered from excessive vibration. Results were obtained using MPI (Message Passing Interface) on a high performance computing (HPC) cluster, evolving 1000 different LSTM cell structures using 168 cores over 4 days. The new evolved LSTM cells showed an improvement of 1.35%, reducing prediction error from 5.51% to 4.17% when predicting excessive engine vibrations 10 seconds in the future, while at the same time dramatically reducing the number of weights from 21,170 to 11,810.

*Keywords:* Ant Colony Optimization, ACO, Long Short Term Memory Recurrent Neural Network, LSTM, Recurrent Neural Network, RNN, Time Series Prediction, Aviation, Aerospace Engineering, Turbomachinery, Turbine engine vibration, Flight Parameters Prediction
*2010 MSC:* 00-01, 99-00

## 1. Introduction

Aircraft engine vibration is of critical interest to the aviation industry, and accurate predictions of excessive engine vibration have the potential to save time, effort, money as well as human lives in the aviation industry. An aircraft engine, as turbo-machinery, should normally vibrate as it has many dynamic parts. However, it is not supposed to exceed resonance limits as to not destroy the engines [1]. As an example, A. V. Srinivasan [1] describes vibrations generated from engine blades' fluttering. Engine blades are the rotating engine components that have the largest dimensions among the other engine components. When rotating at high speeds, they will withstand high centrifugal forces that would logically give the highest contribution to engine vibrations. Engine vibrations are not that simple to calculate or predict analytically because of the fact that various parameters contribute to their occurrence. This fact is always a problem for aviation performance monitors, especially as engines vary in design, size, operation conditions, service life span, the

aircraft they are mounted on, and many other parameters. Most of these parameters contributions can be translated in some key parameters measured and recorded on the flight data recorder. Nonetheless, vibrations are likely to be a result of a mixture of these contributions, making it very hard to predict the real cause behind the excess in vibrations.

As such, engine vibration is a complex problem depending on an unknown number of parameters interacting over an unknown extended period of time, which poses a challenge in analyzing its causes and triggers. Holistic computation methods represent a promising solution for this problem by letting the computers find relations and anomalies that might lead to the problem through a learning process using time series data from flight data recorders (FDR). Traditional neural networks, however, lack the required capabilities to capture those relations and anomalies as they work on current time series without taking the effect of the previous time instances' parameters on the current or future time instants. Due to this, recurrent neural networks have been developed which utilize memory neurons that retain information from previous passes for use with the current experienced data, giving a chance for the neural network to know which parameter really have higher contributions to the investigated problem.

However, these complicated neural network designs in turn posses their own challenges. Regardless of the difficulty of implementing it to a specific problem, the learning process is the main concern when dealing with such neural networks with a large number of interactive connections. When supervised learning is considered and the back-propagation is implemented to update the weights of the connections of the neural network, vanishing and exploding gradients are very serious obstacles for the successfully training recurrent neural networks. As noted by Hochrieter and Schmidhuber [2], *"Learning to store information over extended period of time intervals via recurrent backpropagation takes a very long time, mostly due to insufficient, decaying error back flow."*

While this drawback hindered the application of such sophisticated neural network designs, RNNs which utilize LSTM memory cells offer a solution for this problem as the memory cells provide forget and remember gates which prevent or lessen vanishing or exploding gradients. LSTM RNNs have been used successfully in many studies on involving time series data [3, 4, 5, 6, 7] and were chosen by this study to examine them as a solution to predicting aircraft engine vibration.

For many years, neural networks have strongly proven their prediction potential and applied in different many areas [8, 9, 10]. Various strategies and techniques exist to automatically generate the structure of neural networks and the most common used ones are based on evolutionary algorithms. These neuroevolution strategies examine the use of learning and evolution as concepts to improve the performance of neural networks. In traditional neuroevolution [11], an evolutionary algorithm is used to train a neural networks' connection weights with a fixed structure, but significant benefit has been demonstrated in using these techniques to both optimize and evolve connections and topologies [12, 13, 14], as weights are not the only key parameter for best performance of neural networks [15]. These strategies are of particular interest as determining the optimal structure for a neural network is still an open question. This particular work focuses on evolving the structure of LSTM neurons with an ant colony optimization [16] based algorithm.

### 1.1. Previous work

This study's ultimate goal is to explore the utilization of LSTM RNNs to predict future engine vibration in order to be used in a warning system to give indications for the problem before it occurs in order to avoid or mitigate it. An initial work examined building viable Recurrent Neural Networks (RNN) using Long Short Term Memory (LSTM) neurons to predict aircraft engine vibrations [17]. To achieve this, three different LSTM RNNs architectures were examined to find which would provide better results. The three architectures varied in both complexity and depth of layers. The different networks were trained on time series flight data records obtained from a regional airline containing flights that suffered from excessive vibration. The structure of the LSTM RNNs used in this study is shown in Figure 6. After selecting an initial set of 15 relevant parameters, these LSTM RNNs were able to predict vibration values for 1, 5, 10, and 20 seconds in the future, with 2.84% 3.3%, 5.51% and 10.19% mean absolute error, respectively.

### 1.2. Ant Colony Optimization

Ant colony optimization (ACO) is a metaheuristic used to find approximate solutions to many combinatorial problems. It belongs to a family of bio inspired metaheuristics. ACO is a distributed approach using

agents called artificial ants. These artificial ants resemble biological ants, in that each ant is independent and communicates with other members of the colony through a chemical called pheromone. Ants randomly explore areas, however they tend to follow paths with pheromones, and upon finding food they mark their return path with more pheromone. Pheromones decay over time, and paths with the most pheromone represent the most promising paths to food. The first algorithm of this type (the "Ant System" [18]) was designed for the traveling salesman problem, but failed to produce competitive results. However, subsequent research has shown this algorithm to be effective on this problem [19, 20, 21] and interest for the metaphor has launched many algorithms inspired by it in various fields, including continuous optimization problems [22, 23, 24, 25, 26, 27], and even training neural networks [28, 29, 30, 31].

### 1.3. Study's Contribution

This work improves the performance of the LSTM network architectures proposed in the previous work by optimizing LSTM cell structure. Since the contributions of the input parameters are not of the same magnitude to the problem (vibration), the weights of the neural network are adjusted through backpropagation. However, a fully connected neural network can pose additional training complexity and unwanted noise through some of the connections coming out of some of the inputs. Therefore, there is a need for a way to examine these connections and try to eliminate those which contribute to the prediction error. ACO has been chosen mainly because it has proven its effectiveness in evolving RNNs [16] for time series data prediction. The results have demonstrated that the evolved LSTM architecture increased the best previous results' performance by 1.35% in predicting vibration 10 seconds in the future, while at the same time only requiring nearly half of the connections (the number of weights was reduced from 21,170 to 11,810). The prediction accuracy was improved from 94.49% to 95.83% based on mean absolute error calculations.

## 2. Related Work

### 2.1. Turbine Engine Vibration

According to A. V. Srinivasan [1]: *"The most common types of vibration problems that concern the designer of jet engines include (a) resonant vibration occurring at an integral order, i.e. multiple of rotation speed, and (b) flutter, an aeroelastic instability occurring generally as a nonintegral order vibration, having the potential to escalate, underline{unless checked by any means available to the operator}, into larger and larger stresses resulting in serious damage to the machine. The associated failures of engine blades are referred to as high cycle fatigue failures"*. The *means available to the operator* in practical aviation operations are mainly the implementation of manufacturers' maintenance program which relies on reliability observations.

Any aircraft maintenance program has four objectives: *i)* guaranteeing the inherit safety and reliability levels of the systems and subsystems, *ii)* restore those levels to their original levels if deviations occur, *iii)* gather information about design enhancement for systems and subsystems that showed deficiency in there expected reliability, and *iv)* accomplish these targets at the the most overall cost efficiency possible. To achieve these goals a maintenance program consists of checks performed over specific intervals to perform on-condition checks which might be visual inspections, test runs, or non-destructive tests performed on the systems' or subsystems components. In addition to that, and as mentioned earlier, preventive maintenance is also part of the maintenance program where parts are replaced or overhauled based on reliability observations [32].

There are also other methods to mitigate the risk coming from excessive engine vibration using statistical and holistic computation methods. This is accomplished by monitoring the engine performance through its flight data history, which is logged as time series data, in order to forecast the excessive vibration occurrence. However, since these methods are not exact, reasonable safety factors should be considered. As discussed in the following sections, machine learning and artificial intelligence becomes the heart of such methods.

### 2.2. Time Series Data Prediction

From a statistical point of view, the main goal of prediction is to provide vital information for decision makers, economists, planners optimizers, industrialists and critical systems operators. There are two sides for prediction: the qualitative side and the quantitative side. The qualitative side utilizes methods known as the judgmental or subjective prediction methods which covers methods relaying on intuition, judgement or opinions of some kind of a referee as consumers, experts and/or supporting information. Qualitative

methods are considered in cases when past data is not available. On the other hand, quantitative methods include univariate and multivariate methods. For many study cases related to different scientific and real life problems, the time series data are available on several dependent variables, and in such cases multivariate prediction methods are used [33].

The models in the time series predictions realm mainly falls in two categories: *a)* statistical prediction models which include, *e.g.*, the autoregressive (AR) model, the moving average (MA) model, and hybrid models that derive from them such as autoregressive moving average (ARMA), autoregressive integrated moving average (ARIMA), seasonal ARMIA (SARIMA) [34], vector autoregressive (VAR) models and *b)* artificial neural networks (ANN) prediction models. While successful studies have managed to achieve good results by using such methods and even reported results that out perform neural networks [35], the main draw back of the statistical methods is that they generally can not be applied to non-linear systems [34]. On the other hand, ANNs have shown good performance with such systems. There are also a third category which are the hybrid models. These models (hybrid) offer the benefits of both statistical and ANN models. Consequently, as the studied system involving prediction of aircraft engine vibration is complicated in its nature and is expected to be extremely non-linear, statistical and hybrid models are considered beyond the scope of the study.

### 2.3. Time Series Prediction in Aviation

Some effort has been done using neural networks to classify engine abnormalities without doing analytical computation, *e.g.*, Alexandre Nairac *et al.* [36] have performed research to detect abnormalities in engine vibrations based on recorded data. To achieve that, the work used two modules. One of the modules uses the overall shape of the vibration curve to detect unusual vibration signatures. The second one reports sudden unexpected transitions in the signature curves. Their approach to detect defects is not to introduce examples of faulty engines to the neural network, rather, only examples of healthy engines are introduced to the neural networks in the training phase. This approach was taken to overcome the lack of existence of adequate faulty engine data, which was not enough for training. In this context, the paper introduces the term 'normality' to describe the behavior of normal engines and 'abnormality' to describe the behavior of faulty engines. Using statistical models, the faulty engines detection would be described as 'novelty' detection based on the deviation from the data distribution. The best results this work achieved was the prediction of faulty engines with 84% successful classifications.

David A. Clifton *et al.* [37] presented work for predicting abnormalities in engine vibration based on statistical analysis of vibration signatures. The paper presents two modes of prediction. One is ground-based (off-line), where prediction is done by run-by-run analysis to predict abnormalities based on previous engine runs. The success in this approach was predicting abnormalities two flights ahead. The other mode is a flight based-mode (online) in which detection is done either by sending reduced data to the ground-base or processing it onboard the aircraft. The paper mentions that they could successfully predict vibration events 2.5 hours in the future. However, this prediction is done after half an hour of flight data collection, which might be a critical time as well, as excess vibration may occur during this data collection time. The paper did not mention how much data was required to have a sound prediction.

### 2.3.1. RNN for Predicting Flight Parameters

Having an advantage over standard FFNNs[1], RNNs can deal with sequential input data, using their internal memory to process sequences of inputs and use previously stored information to aid in future predictions. This is done by feedback connections or by looping between neurons, which allows them to be of predicting more complex data [38].

This presented work is in part inspired by previous work on predicting flight parameters [39, 16]. Which first utilized evolutionary algorithms such as particle swarm optimization [40, 41] and differential evolution [42] to optimize the weights of the network [39], and then an ant colony optimization based algorithm to evolve different recurrent neural network structures [16]. The neural networks evolved with ant colony optimization predicted airspeed, altitude and, pitch with a 63%, 97% and 120% improvement respectively over

---

[1]Feed Forward Neural Networks

the previously best published results. The research used recurrent neural networks and applied an ant-colony optimization (ACO) algorithm [43, 44, 45], an optimization technique used in the beginning on discrete problems, mainly on the Traveling Salesman Problem [46]. Later it was used in continuous optimization problems [22, 23, 24, 25, 26, 27], including training neural networks [28, 29, 30, 31].

### 2.3.2. LSTM RNN

LSTM RNNs were first introduced by S. Hochrieter & J. Schmidhuber [4]. While the work by T. Desell *et al.* utilized non-gradient based evolutionary algorithms to optimize RNN weights, LSTM neurons provide a solution for the exploding/vanishing gradients problem by utilizing various gates, which allow backpropagation to be used in large RNNs (S. Hochrieter in 1991). This work has paved the way for many interesting projects.

Later, J. Schmidhuber *et al.* [47] emphasized the forget gate in the LSTM RNNs. The paper mentions that *"We identify a weakness of LSTM networks processing continual input streams that are not a priori segmented into subsequences with explicitly marked ends at which the network's internal state could be reset. Without resets, the state may grow indefinitely and eventually cause the network to break down. Our remedy is a novel, adaptive forget gate that enables an LSTM cell to learn to reset itself at appropriate times, thus releasing internal resources. We review illustrative benchmark problems on which standard LSTM outperforms other RNN algorithms. All algorithms (including LSTM) fail to solve continual versions of these problems. LSTM with forget gates, however, easily solves them, and in an elegant way."* However, Felix A. Gers *et al.* [48] suggest that *"LSTM RNNs does not carry over to certain simpler time series prediction tasks solvable by time window approaches"*. The paper suggests to use LSTM when *"simpler traditional approaches fails"*.

LSTM RNNs have been used with strong performance in image recognition [49], audio visual emotion recognition [50], music composition [51] and other areas. Regarding time series prediction, for example, LSTM RNNs have been used for stock market forecasting [3] and forex market forecasting [7]. Also forecasting wind speeds [4, 5] for wind energy mills, and even predicting diagnoses for patients based on health records [6].

### 2.4. Evolutionary Optimization Methods

Several methods for evolving topologies along with weights have been searched and deployed. In [52, 52], NeuroEvolution of Augmenting Topologies (NEAT) has been developed. It is a genetic algorithm that evolves increasingly complex neural network topologies, while at the same time evolving the connection weights. Genes are tracked using historical markings with innovation numbers to perform crossover among different structures and enable efficient recombination. Innovation is protected through speciation and the population initially starts small without hidden layers and gradually grows through generations [53, 54, 55]. Experimentations have demonstrated that NEAT presents an efficient way for evolving neural networks for weights and topologies in parallel or separately. Its power resides in its ability to combine all the four main aspects discussed above and expand to complex solutions along the generation process. However NEAT still has some limitations when it comes evolving neural networks with weights or LSTM cells for time series prediction tasks as it has been claimed in [16].

## 3. Methodology

This work utilizes the ACO method developed by Desell *et al.* to evolve the structure of LSTM cells, in part due to strong previous results and because it allows any method to be used to determine the optimal weights of connections. This is particularly important as it allows backpropagation to be used on a large scale LSTM RNN, which is significantly more efficient than the non-gradient based evolutionary algorithms used in previous work.

### 3.1. Experimental Data

The flight data used consists of 76 different parameters recorded on the aircraft Flight Data Recorder (FDR), inclusive of the engine vibration parameters. During the data processing phase of the project, two efforts were done to identify the parameters that most contributed to the engine vibration.
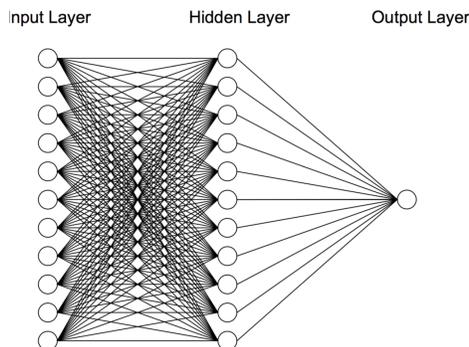
Figure 1: A one layer feed forward neural network.

### 3.1.1. Data Correlation Parameter Selection

Primarily, cross-correlation analysis [56] was exercised to find the potential parameters that highly contribute to vibration. Every parameter from each flight was cross-correlated to vibration then plotted to pick the highest correlated parameters. Cross correlation was calculated using the following Equation:

$$Cross\_Correlation = \sum_{a=-\infty}^{\infty} x[a] \cdot vib[a] \tag{1}$$

Where the highest correlation was determined by calculating the area under the plotted curve for the normalized data. The top correlated parameters to vibration were:

1. Right InBoard Spoiler
2. Right OutBoard Spoiler
3. Left InBoard Spoiler
4. Left OutBoard Spoiler
5. Static Air Temperature
6. Pitch 2
7. Pitch
8. Slat Configuration
9. Main Landing Gear Lock Down Sensors
10. Flap Configuration

A one layer feed forward neural network was built as shown in Figure 1 to predict vibration given other parameters within the same second. However, the results were poor, with significant noise in the predictions. This imposed a question about the quality of the chosen parameters using this method and due to this, another method of parameter-selection was sought. A potential cause for such misleading cross-correlation chosen parameters was that some flight configuration parameters like spoilers/slats/flaps positions, pitch angle and main-landing-gear position do not change but few times during the flight, which can translate into high correlation with the vibration.

### 3.1.2. Aerodynamics/Turbo-machinery Parameter Selection

A second subset of the FDR parameters were then chosen based on the likelihood of their contribution to the vibration based on aerodynamics/turbo-machinery expert knowledge. Again, a one layer feed forward neural network with a structure similar to the one shown in Figure 1, except the number of input and hidden nodes was equal to the number of chosen parameters, was applied and these results were encouraging enough to use these parameters for predicting vibration in future.

Some parameters, such as Inlet Guide Vans Configuration, Fuel Flow, Spoilers Configuration (this was preliminarily considered because of the special position of the engine mount), High Pressure Valve Configuration and Static Air Temperature were excluded because it was found that they generated more noise than positively contributing to the vibration prediction.
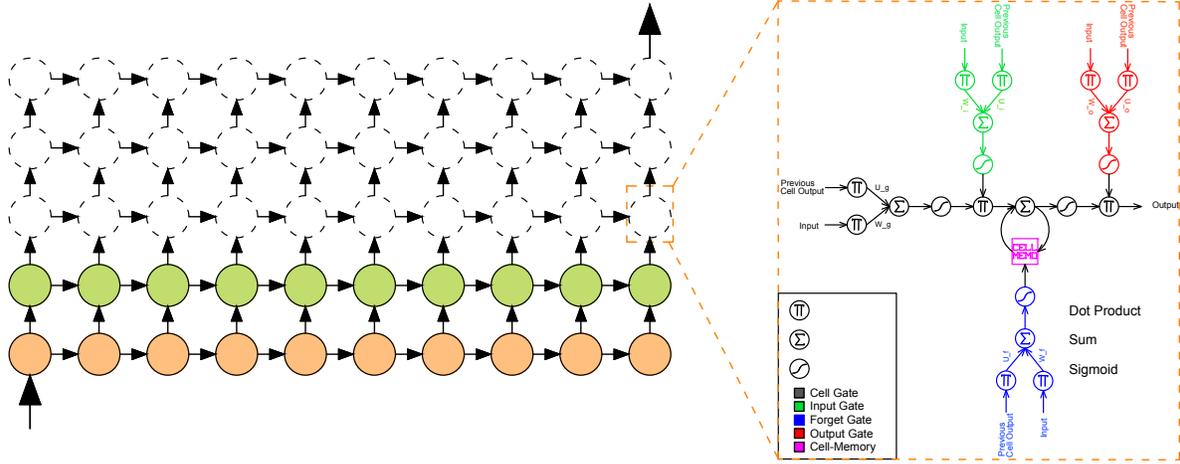
6

Figure 2: A generic overview of the design of a LSTM RNN.

The final chosen parameters were:

1. Altitude [ALT]
2. Angle of Attack [AOA]
3. Bleed Pressure [BPRS]
4. Turbine Inlet Temperature [TIT]
5. Mach Number [M]
6. Primary Rotor/Shaft Rotation Speed [N1]
7. Secondary Rotor/Shaft Rotation Speed [N2]
8. Engine Oil pressure [EOP]
9. Engine Oil Quantity [EOQ]
10. Engine Oil Temperature [EOT]
11. Aircraft Roll [Roll]
12. Total Air Temperature [TAT]
13. Wind Direction [WDir]
14. Wind Speed [WSpd]
15. Engine Vibration [Vib]

*3.2. Recurrent Neural Network Design*

Three LSTM RNN architectures were designed to predict engine vibration 5 seconds, 10 seconds, and 20 seconds in the future. Each of the 15 selected FDR parameters is represented by a node in the inputs of the neural network and an additional node is used for a bias. Each neural network in the three designs consists of LSTM cells that receive both an initial input of flight data at some time in the past or the output from a cell in the lower layer, and the output of the previous cell in the same layer, as inputs (see Figure 2). Each cell has three gates to control the flow of information through the cell and accordingly, the output of the cell. Each cell also has a cell-memory which is the core of the LSTM RNN design. The cell-memory allows the flow of information from the previous states into the current predictions.

The gates that control the flow are shown in Figure 3. They are: *i)* the *input gate*, which controls how much information will flow from the inputs of the cell, *ii)* the *forget gate*, which controls how much information will flow from the cell-memory, and *iii)* the *output gate*, which controls how much information will flow out of the cell. This design allows the network to learn not only about the target values, but also about how to tune its controls to reach the target values.

All the utilized architectures follow the common LSTM RNN designs shown in Figure 2 and 3. However, there are two variations of this common design used in the utilized architectures, shown in Figures 4 and 5,
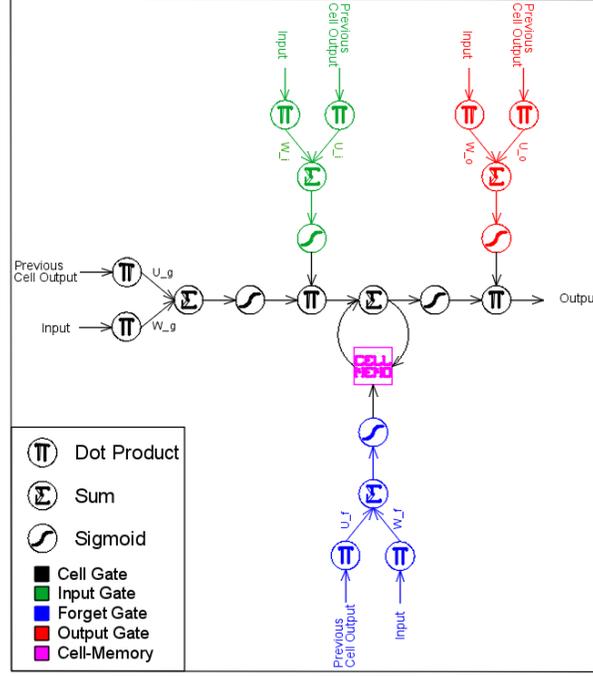
7

Figure 3: LSTM cell design

with the difference being the number of inputs from the previous cell. Cells that take an initial number of inputs and output the same number of outputs are denoted by *'M1'* cells. As input nodes are needed to be reduced through the neural network, the design of the cells are different. Cells which perform a reduction on the inputs are denoted by *'M2'* cells.

### 3.2.1. LSTM RNN Forward Propagation Equations

The equations used in the forward propagation through the neural network are:

$$i_t = Sigmoid(w_i \bullet x_t + u_i \bullet a_{t-1} + bais_i) \tag{2}$$

$$f_t = Sigmoid(w_f \bullet x_t + u_f \bullet a_{t-1} + bais_f) \tag{3}$$

$$o_t = Sigmoid(w_o \bullet x_t + u_o \bullet a_{t-1} + bais_o) \tag{4}$$

$$g_t = Sigmoid(w_g \bullet x_t + u_g \bullet a_{t-1} + bais_g) \tag{5}$$

$$c_t = f_t \bullet c_{t-1} + i_t \bullet g_t \tag{6}$$

$$a_t = o_t \bullet Sigmoid(c_t) \tag{7}$$

where (see Figure 3):

$i_t$: input-gate output
$f_t$: forget-gate output
$o_t$: output-gate output
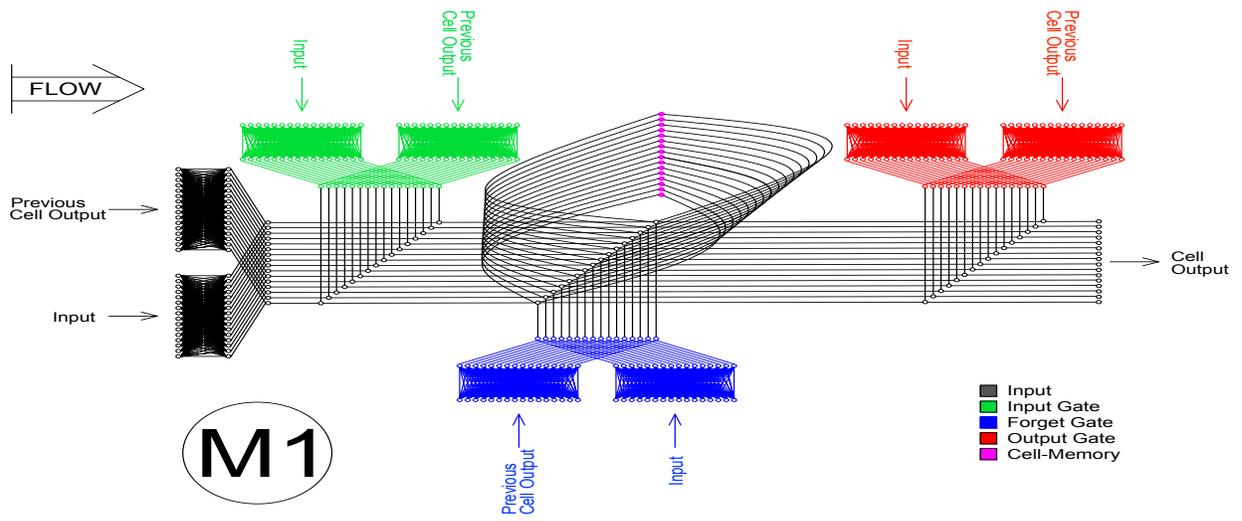$g_t$: input's sigmoid
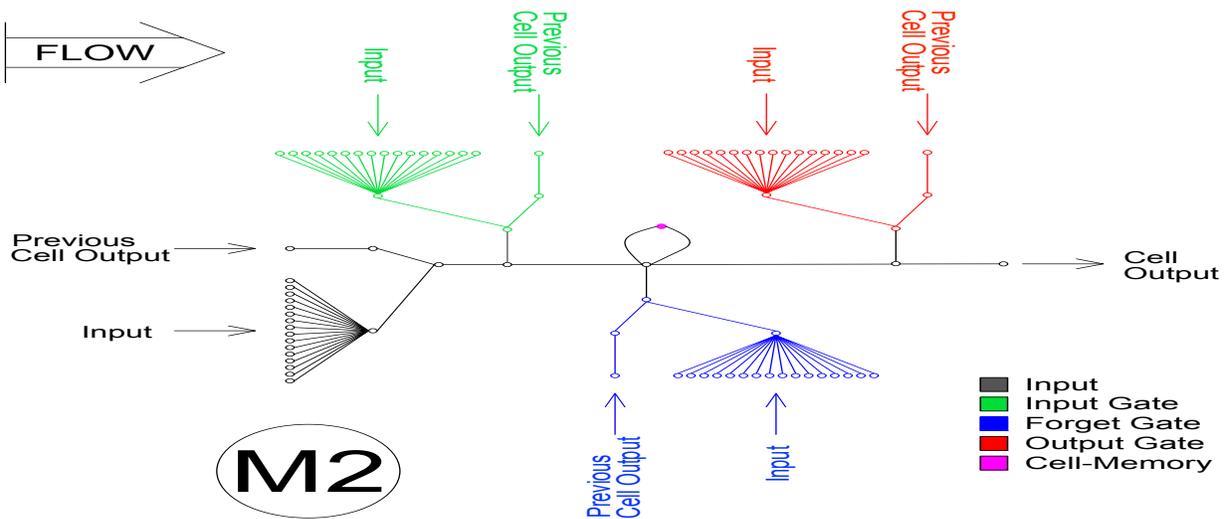$c_t$: cell-memory output

8

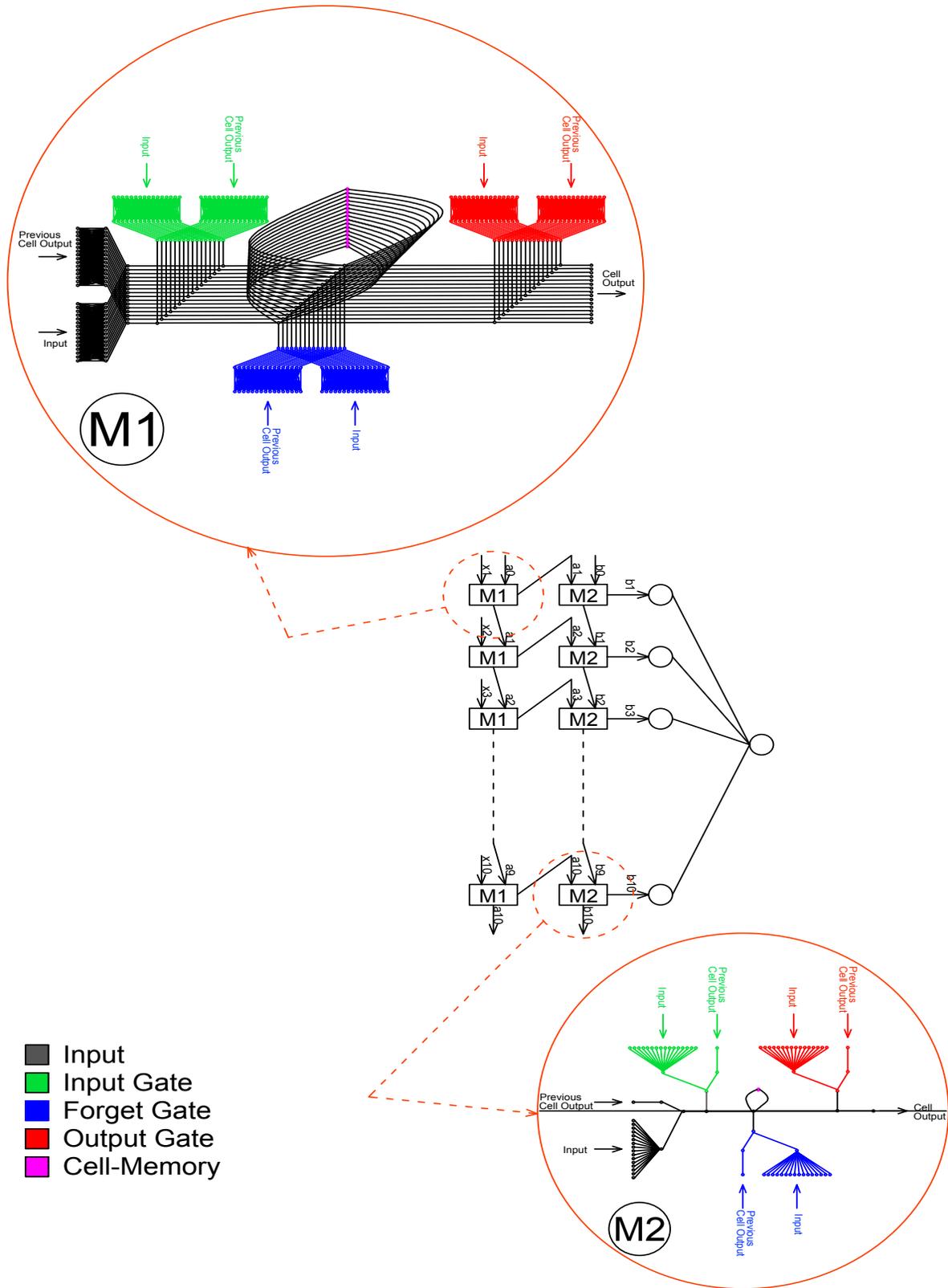Figure 4: Level 1 LSTM cell design



Figure 5: Level 2 LSTM cell design

Figure 6: Neural network structure

$w_i$: weights associated with input and input-gate

$u_i$: weights associated with previous output and input-gate

$w_f$: weights associated with input and forget-gate

$u_f$: weights associated with previous output and forget-gate

$w_o$: weights associated with input and output-gate

$u_o$: weights associated with previous output and the output-gate

$w_g$: weights associated with the cell input

$u_g$: weights associated with previous output and the cell input

and the formula of the sigmoid function is:

$$Sigmoid(\alpha) = \frac{1}{1 + e^{-\alpha}} \tag{8}$$

### 3.3. LSTM RNN Architectures

The three architectures are as follows, with the dimensions of the weights of these architectures shown in Table 1 and the total number of weights shown in Table 2.

**Architecture I**

As shown in Figure 7a, the first level of the architecture takes inputs from ten time series (the current time instant and the past nine). It then feeds the second level of the neural network with the output of the first level. The output of the first level of the neural network is considered the first hidden layer. The second level of the neural network then reduces the number of nodes fed to it from 16 nodes (15 input nodes + bias) per cell to only one node per cell. The output of the second level of the neural network is considered the second hidden layer. Finally, the output of the second level of the neural network would be only 10 nodes, a node from each cell. These nodes are fed to a final neuron in the third level to compute the output of the whole network.

The dimensions of the weights matrices and vectors of this architecture are shown in Table 1. The total number of weights are is shown in Table 2. Figures 8 and 9 provide an overview of architecture I, as it has a large number of connections (21,170). Figure 8 shows the overall design of how the LSTM cells are connected, and then Figure 9 displays all the connections within a single time step of the full LSTM RNN. As a whole, there are 10 different instances of Figure 8, each connected as specified in Figure 9.
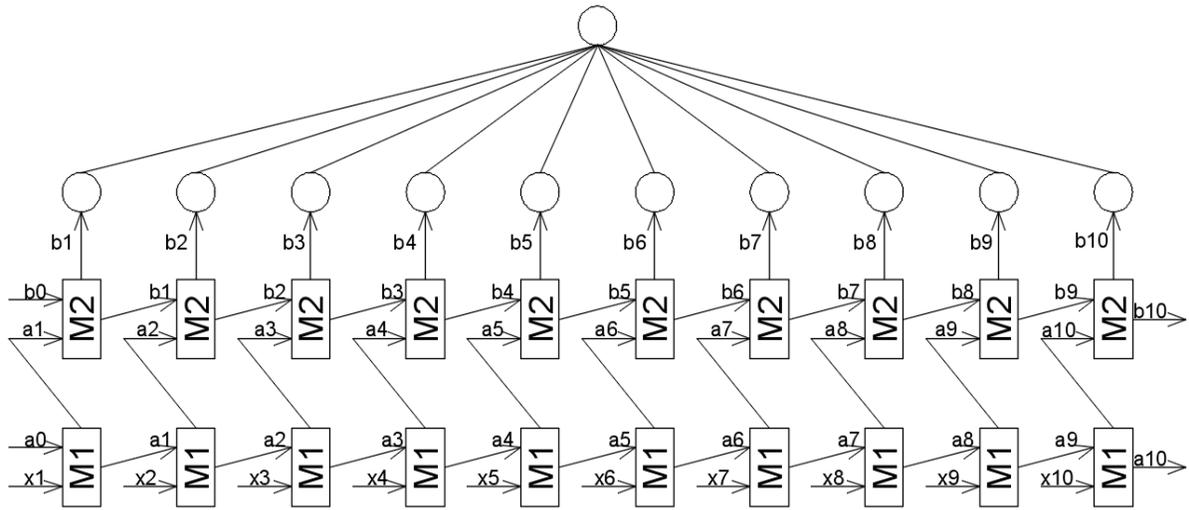
**Architecture II**

As shown in Figure 7b, this architecture is almost the same as the previous one except that it does not have the third level. Instead, the output of the second level is averaged to compute the output of the whole network.

The dimensions of the weights matrices and vectors of this architecture are shown in Table 1. The total number of weights are shown in Table 2.

**Architecture III**

Figure 7c presents a deeper neural network architecture. In this design, the neural network takes inputs from twenty time series (the current time instant and the past nineteen) as the first level. It feeds the second level of the neural network with the output from the first level. The second level does the same procedure as first level giving a chance for more abstract decision making. The output of the second level of the neural network is considered the first hidden layer and the output of the second level is considered the second hidden layer. The third level of the neural network then reduces the number of nodes fed to it from 16 nodes (15 input nodes + bias) per cell to only one node per cell. The output of the third level of the neural network is considered the third hidden layer. Finally, the output of the third level of the neural network is twenty nodes, a node from each cell. These nodes are fed to a final neuron in the fourth level to compute the output of the whole network.

The Dimensions of the weights matrices and vectors of this architecture are shown in Table 1. The total number of weights are shown in Table 2.
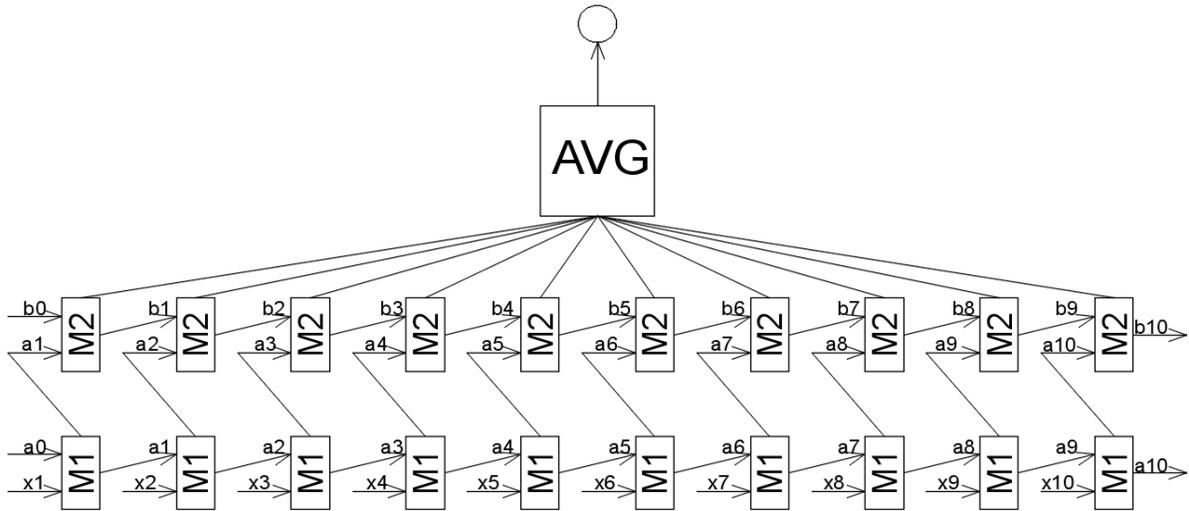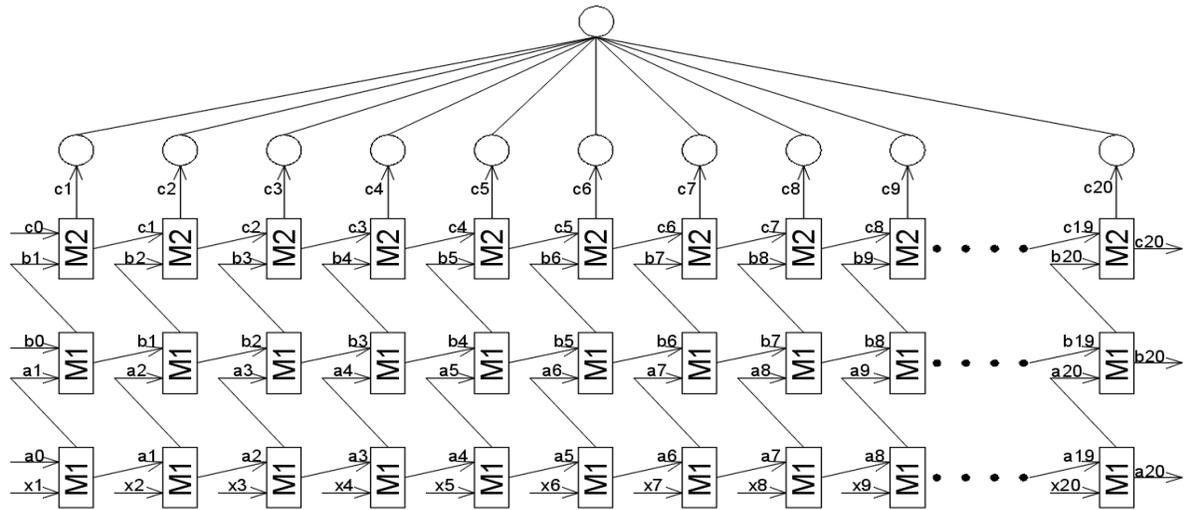
(a) Architecture I

(b) Architecture II

(c) Architecture III

Figure 7: Used LSTM RNNs Architectures

Figure 8: One time step of Architecture I

Table 1: Architectures Weights-Matrices Dimensions

**Architecture I**

| | $w_i$ | $u_i$ | $w_f$ | $u_f$ | $w_o$ | $u_o$ | $w_g$ | $u_g$ |
|---|---|---|---|---|---|---|---|---|
| **Level 1** | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ |
| **Level 2** | $16\times1$ | $1\times1$ | $16\times1$ | $1\times1$ | $16\times1$ | $1\times1$ | $16\times1$ | $1\times1$ |
| **Level 3** | | | | $16\times1$ | | | | |

**Architecture II**

| | $w_i$ | $u_i$ | $w_f$ | $u_f$ | $w_o$ | $u_o$ | $w_g$ | $u_g$ |
|---|---|---|---|---|---|---|---|---|
| **Level 1** | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ |
| **Level 2** | $16\times1$ | $1\times1$ | $16\times1$ | $1\times1$ | $16\times1$ | $1\times1$ | $16\times1$ | $1\times1$ |

**Architecture III**

| | $w_i$ | $u_i$ | $w_f$ | $u_f$ | $w_o$ | $u_o$ | $w_g$ | $u_g$ |
|---|---|---|---|---|---|---|---|---|
| **Level 1** | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ |
| **Level 2** | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ | $16\times16$ |
| **Level 3** | $16\times1$ | $1\times1$ | $16\times1$ | $1\times1$ | $16\times1$ | $1\times1$ | $16\times1$ | $1\times1$ |
| **Level 4** | | | | $16\times1$ | | | | |

13

Figure 9: One time step of the Architecure I Full Structure

Table 2: Architectures Weights Matrices' Total Elements

| Architecture I | Architecture II | Architecture III |
|:---:|:---:|:---:|
| 21,170 | 21,160 | 83,290 |

### 3.4. Forward Propagation

The following is a general description for the forward propagation path. This example uses Architecture I as an example but similar steps are taken in the other architectures with minor changes apparent in their diagrams. With Figure 7a presenting an overview of the structure of the whole network and considering Figure 4 as an overview of the structure of the cells in *Level 1* (M1) and Figure 5 as an overview of the structure of the cells in *Level 2* (M2) – the input at each iteration consists of 10 seconds of time series data of the 15 input parameters and 1 bias (*Input* in Figure 4) in one vector ($x_t$ in Figure 7a) and the output of the previous cell (*Previous Cell Output* in Figure 4) in another vector ($a_{t-1}$ in Figure 7a). Each second of time series input is fed to the corresponding cell (*i.e.*, the first seconds' 15 parameters and 1 bias are fed to first cell, the second seconds' 15 parameters and 1 bias are fed to second cell, ...) into the *cell gate* (shown in black color), *input gate* (shown in green color), *forget gate* (shown in blue color) and the *output gate* (shown in red color). If the gates (*input gate*, *forget gate* and, *output gate*) are seen as valves that control how much of the data flow through it, the outputs of these gates ($i_t$, $f_t$ and, $o_t$) are considered as how much these valves are opened or closed.

First, at the *cell gate*, $x_t$ is dot multiplied by its weights matrix $w_g$ and $a_{t-1}$ is dot multiplied by its weights matrix $u_g$. The output vectors are summed and an activation function is applied to it as in Equation 5. The output is called $g_t$.

Second, at the *input gate*, $x_t$ is dot multiplied by its weights matrix $w_i$ and $a_{t-1}$ is dot multiplied by its weights matrix $u_i$. The output vectors are summed and an activation function is applied to it as in Equation 2. The output is called $i_t$.

Third, at the *forget gate*, $x_t$ is dot multiplied by its weights matrix $w_f$ and $a_{t-1}$ is dot multiplied by its weights matrix $u_f$. The output vectors are summed and an activation function is applied to it as in Equation 3. It controls how much of the *cell memory* Figure 7a (saved from previous time-step) should pass. The output is called $f_t$.

Fourth, at the *output gate*, $x_t$ is dot multiplied by its weights matrix $w_o$ and $a_{t-1}$ is dot multiplied by its weights matrix $u_o$. The output vectors are summed and an activation function is applied to it as in Equation 4. The output is called $o_t$.

Fifth, the contribution of the cell input *Input* $g_t$ and *cell memory* $c_{t-1}$ is decided in Equation 6 by dot multiplying them by $f_t$ and $i_t$ respectively. The output of this step is the new *cell memory* $c_t$.

Sixth, cell output is also regulated by the output gate (valve). This is done by applying the sigmoid function to the *cell memory* $c_t$ and dot multiplying it by $o_t$ as shown in Equation 7. The output of this step is the final output of the cell at the current time-step $a_t$. $a_t$ is fed to the next cell in the same level and also fed to the cell in the above level as an *Input* $a_t$.

The same procedure is applied at *Level 2* but with different weight vectors and different dimensions. Weights at *Level 2* have smaller dimensions to reduce their input detentions from vectors with 16 dimensions to vectors with one dimension. The output from *Level 2* is a one dimensional vector from each cell of the 10 cells in *Level 2*. These vectors are fed as one 10 dimensional vector to a simple neuron collection shown in Figure 7a at *Level 3* to be dot multiplied by a weight vector to reduce the vector to a single scalar value: the final output of the network at the time-step.

## 4. Evolving LSTM RNN Cells using Ant Colony Optimization

Although the results from architecture I are promising, there is still room for further optimization in that the network may have excessive connections which confound accurate predictions and that the structure could be further optimized. A particular concern was that some connections could cause noise in the obtained results and ultimately would drift the results from their most optimum values, as this had been shown in the initial one layer feed forward neural networks with certain input parameters. The goal of using the ant
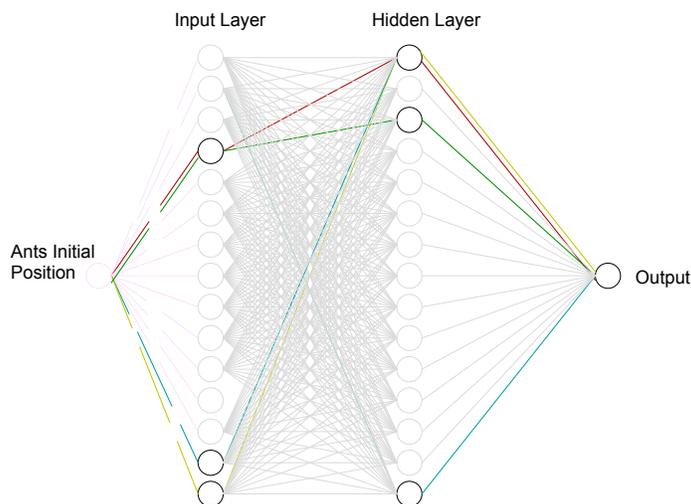
Figure 10:   Schematic of Neural Network Structure after AOC

colony optimization strategy is to evolve the structure of the LSTM cells, encouraging more diverse networks and selecting the topologies that give the best performance.

The ACO algorithm operates on the fully connected inputs to the M1 and M2 cells, as shown in Figures 4 and 5. Each M1 cell has eight 16x16 input gates, four of which take the input from the previous cell in the same layer, and four of which take the input from the time series or the cell in the lower layer. Each M2 cell has eight 8x1 input gates, four of which receive input from the previous cell in the same layer and four from the cell in the lower layer.

The algorithm begins with a fully connected gate that will be used by the ants each time to generate new paths for new network designs. Paths are selected by the ants based on pheromones – each connection in the network has a pheromone value that determine its probability to be chosen as a path. Given a number of ants, each one will select one path from the fully connected network. All the paths selected from all the ants are then collected, duplicated paths are removed and a design network is generated based on the new cell topology. Figure 10 shows an example on an M1 cell, assuming four ants choosing their paths on an input gate to an M1 cell, which generates a subgraph from the potentially fully connected input gate. The same ACO generated topology is used for each of these 8 input gates. Figure 17a provides an example of the best found ACO optimized M1 cell.

In detail, the paths generated by ACO are used in the connections between the *"Input"* and the hidden layer neurons that follow it, and the *"Previous Cell Output"* and the hidden layer neurons that follow it. The connections between the *"Input"* and the hidden layer neurons that follow it are shown in the first level cells' (Figure 4 "M1") in BLACK color, GREEN color, BLUE color, and RED color at the gates of the cell. Once a hidden node in first level cell is reached by an ant, the connection between this node and the output node shown in second level cells' Figure 5 "M2" in BLACK color, GREEN color, BLUE color, and RED color, will automatically be part of the evolved mesh because the ant will not have any other option to reach the output node except through that single connection.

The same generated mesh is used at all the gates: Main Gate, Input Gate, Forget Gate, and Output Gate at the "M1"cells and "M2" cells at all the time-steps in the LSTM RNN Architecture I as shown in Figure 7a. In other words, regardless the LSTM RNN time-step, whenever there is a transition without data reduction: the first set of connections in the generated mesh is used, and whenever there is a transition with data reduction: the second set of connections in the generated mesh is used.

### 4.1. Distributed ACO Optimization

Evolving large LSTM RNNs is a computationally expensive process. Even training a single LSTM RNN is extremely time consuming (approximately 8.5-9 hours to train one architecture), and applying the ACO

algorithm requires running the training process on each evolved topology. This significantly raises the computational requirements in time and resources necessary to process and evolve better LSTM networks. For that reason, the ant colony algorithm was parallelized using the message passing interface (MPI) for Python [57] to allow for it to be run utilizing high performance computing resources.

The distributed algorithm utilizes an asynchronous master worker approach, which has been shown to provide performance and scalability over iterative approaches in evolutionary algorithms [58, 59]. This approach provides an additional benefit in that it is automatically load balanced – workers request and receive new LSTM RNNs when they have completed training previous ones, without blocking on results from other workers. The master process can generate a new LSTM RNN to be trained from whatever is currently present in its population.

The algorithm is defined in Algorithm 1. In detail, the algorithm beings with the master process generating an initial set of network designs randomly (given a user defined number of ants), and sending these to the worker processes. When the worker receives a network design, it creates an LSTM RNN architecture by creating the LSTM cells with the according input gates and cell memory. The generated structure is then trained on different flight data records using the backpropagation algorithm and the resulting fitness (test error) is evaluated and sent back along with the LSTM cell paths to the master process.

The master process then compares the fitness of the evaluated network to the other results in the population, inserts it into the population, and will reward the paths of the best performing networks by increasing the pheromones by 15% of their original value if it was found that the result was better than the best in the population. However, the pheromones values are not allowed to exceed a fixed threshold of 20. The networks that did not out perform the best in the population are not penalized by reducing the pheromones along their paths.

## 5. Implementation

### 5.1. Programming Language

Python's Theano Library [60] was used to implement the neural networks. It was chosen due to four major advantages: *i)* it will compile the most, if not all, of functions coded using it to C and CUDA providing fast performance, *ii)* it will perform the weights updates for backpropagation with minimal overhead, *iii)* Theano can compute the gradients of the error (cost function output) with respect to the weights, saving significant effort and time needed to manually derive the gradients, coding and debugging them (which is particularly challenging in regards to LSTM neurons), and finally, *iv)* it can utilize GPUs for further increased performance.

### 5.2. Data Processing

The flight data parameters used were normalized between 0 and 1. The sigmoid function was used as an activation function over all the gates and inputs/outputs. The ArcTan activation function was tested on the data, however it gave distorted results and sigmoid function provided significantly better performance.

### 5.3. Machine Specifications

The algorithm was implemented in Python using MPI for Python [57] and was run on the University of North Dakota's high performance computing cluster. The cluster is running the Red Hat Enterprise Linux (RHEL) 7.2 operating system with 31 nodes, each with 8 cores for 248 in total, 64GBs RAM per node for a total 1948 GB, and it is using InfiniBand 10 gigabit (GB) for interconnect. The same number of epochs (575) were used before to train the LSTM Network as in previous work for comparison purposes.

### 5.4. Using GPUs for LSTM RNN Training

The neural networks' weight matrices for a LSTM cell are repeated at a given time-step at a given layer. Thus, the computational cost increases if the output if these gates is computed separately, one gate at a time, as the data input/output consumes CPU cycles. This case is also obvious if a GPU is utilized for high performance computing as the cost of sending data forward and backward between the CPU (host) and GPU (device). For that, the input of a cell at a given layer is dot multiplied by a matrix that holds all of the gates

**Algorithm 1** Ant Colony Algorithm

1: **global variables**
2:     ▷ A user specified number of ants.
3:     $N\_ANTS$
4:     ▷ A user specified maximum number of pheromones.
5:     $N\_PHEROMONES$
6:     ▷ The number of input parameters (15) + 1 for bias.
7:     $N\_INPUTS \leftarrow 16$
8: **end global variables**
9: **function** GENERATE_PATHS
10:     $paths = $ **new** $Paths$
11:     ▷ path arrays all initialized to 0. 1 indicates a connection.
12:     $paths.input = $ **array**$[16]$
13:     $paths.m1 = $ **array**$[16][16]$
14:     $paths.m2 = $ **array**$[16]$
15:     **for** $ant \leftarrow 1 \dots n\_ants$ **do**
16:         ▷ select input path probabilistically according to pheromones
17:         $pheromone\_sum \leftarrow$ **sum**$(pheromones.input)$
18:         $r \leftarrow$ **uniform_random**$(0, pheromone\_sum - 1)$
19:         $input\_path \leftarrow 0$
20:         **while** $r > 0$ **do**:
21:             **if** $r < pheromones.input[input\_path]$ **then**
22:                 $paths.input\_paths[input\_path] \leftarrow 1$
23:                 **break**
24:             **else**
25:                 $r \leftarrow r - pheromones.input[input\_path]$
26:                 $input\_path \leftarrow input\_path + 1$
27:         ▷ select hidden path probabilistically according to pheromones
28:         $pheromone\_sum \leftarrow$ **sum**$(pheromones.m1[input\_path])$
29:         $r \leftarrow$ **uniform_random**$(0, pheromone\_sum - 1)$
30:         $hidden\_path \leftarrow 0$
31:         **while** $r > 0$ **do**
32:             **if** $r < pheromones.m1[input\_path][hidden\_path]$ **then**
33:                 $paths.m1\_paths[input\_path][hidden\_path] \leftarrow 1$
34:                 $paths.m2\_paths[hidden\_path] \leftarrow 1$
35:                 **break**
36:             **else**
37:                 $r \leftarrow r - pheromones.m1[input\_path][hidden\_path]$
38:                 $hidden\_path \leftarrow hidden\_path + 1$
        **return** paths
39: **function** UPDATE_PHEROMONES(pheromones, paths)
40:     **for** $i \leftarrow 1 \dots paths.input.length$ **do**
41:         **if** $paths.input[i] = 1$ **then**
42:             $pheromones.input[i] \leftarrow$
            $min(pheromones.input[i] * 1.15, MAX\_PHEROMONE)$
43:     **for** $i \leftarrow 1 \dots paths.m1.length$ **do**
44:         **for** $j \leftarrow 1 \dots paths.m1[i].length$ **do**
45:             **if** $paths.m1[i][j] = 1$ **then**
46:                 $pheromones.m1[i][j] \leftarrow$
            $min(pheromones.m1[i][j] * 1.15, MAX\_PHEROMONE)$
        m
47: **procedure** MASTER
48:     ▷ *pheromones all initialized to 1*
49:     $pheromones = $ **new** $Pheromones$
50:     $pheromones.input \leftarrow array[16]$
51:     $pheromones.m1 \leftarrow array[16][16]$
52:     $pheromones.m2 \leftarrow array[16]$
53:     $population = $ **List**$()$

```
54:     repeat
55:         worker, message ← get_next_message()
56:         if message is request_paths then return generate_paths()
57:         else if message is report_fitness then
58:             fitness, paths ← message.get_arguments()
59:             rank ← population.inorder_insert( {fitness, paths} )
60:             if rank = 0 then
61:                 update_pheromones(pheromones, paths)
62:     until finished
63: procedure WORKER
64:     repeat
65:         paths ← master.request_paths()
66:         fitness ← LSTM_RNN(paths).backpropagate()
67:         master.report_fitness(fitness, paths)
68:     until finished
```

weights concatenated one after the other. Then, the outputs; $g$ Equation 5, $i$ Equation 2, $f$ Equation 3 and, $o$ Equation 4, can be extracted from the dot product output matrix. Equation 9 is an example of combining (concatenating) the weights matrices for the LSTM cells' gates of level one in Architecture I. By this, all weights are transferred between the CPU and the GPU as one data structure, which would theoretically boost the performance.

These measures were followed when using the Theano library for GPU computing to manage the GPU threads, blocks and grids as well as the data transfer between the CPU and GPU. However, the performance was slower when compared to the pure CPU version. For Architecture I as an example, one iteration through the network during the learning process, it took the GPU version more than twenty minutes while it took slightly more than two minutes for the pure CPU version.

A further effort was made to overcome the data transfer penalty between the CPU and the GPU. The whole input data set was sent to the GPU as one data structure to avoid the data transfer through the iterations at every time series in the data and to perform those iterations on the GPU. Unfortunately, this also did not help with the performance. Ultimately, a conclusion was reached that the subject matrices are not large enough to overcome the data transfer overhead. Further study is required to determine if it is possible to achieve good performance for these types of LSTM RNNs on GPUs.

$$
\begin{bmatrix}
w_{g_{1,1}} & w_{g_{1,2}} & w_{g_{1,3}} & \cdots & w_{g_{1,16}} \\
\vdots & \vdots & \vdots & \ddots & \vdots \\
w_{g_{16,1}} & w_{g_{16,2}} & w_{g_{16,3}} & \cdots & w_{g_{16,16}} \\
w_{i_{1,1}} & w_{i_{1,2}} & w_{i_{1,3}} & \cdots & w_{i_{1,16}} \\
\vdots & \vdots & \vdots & \ddots & \vdots \\
w_{i_{16,1}} & w_{i_{16,2}} & w_{i_{16,3}} & \cdots & w_{i_{16,16}} \\
w_{f_{1,1}} & w_{f_{1,2}} & w_{f_{1,3}} & \cdots & w_{f_{1,16}} \\
\vdots & \vdots & \vdots & \ddots & \vdots \\
w_{f_{16,1}} & w_{f_{16,2}} & w_{f_{16,3}} & \cdots & w_{f_{16,16}} \\
w_{o_{1,1}} & w_{o_{1,2}} & w_{g_{1,3}} & \cdots & w_{o_{1,16}} \\
\vdots & \vdots & \vdots & \ddots & \vdots \\
w_{o_{16,1}} & w_{o_{16,2}} & w_{o_{16,3}} & \cdots & w_{o_{16,16}}
\end{bmatrix}
\odot
\begin{bmatrix}
x_{11} \\
x_{12} \\
x_{13} \\
\vdots \\
x_{16}
\end{bmatrix}
=
\begin{bmatrix}
out_{g_1} \\
\vdots \\
out_{g_{16}} \\
out_{i_1} \\
\vdots \\
out_{i_{16}} \\
out_{f_1} \\
\vdots \\
out_{f_{16}} \\
out_{o_1} \\
\vdots \\
out_{o_{16}}
\end{bmatrix}
\tag{9}
$$

## 6. Results

### 6.1. Previous Training Results

Training process results from the original three architectures are shown in Table 3. These results are directly proportional to the testing results as will be shown in the results section. The errors shown are mean squared error.

Table 3: Previous Training Results

|  | Prediction Error | | | |
|  | 1 seconds | 5 seconds | 10 seconds | 20 seconds |
| --- | --- | --- | --- | --- |
| Architecture I | 0.000154 | 0.000398 | 0.000972 | 0.001843 |
| Architecture II | 0.001239 | 0.001516 | 0.001962 | 0.002870 |
| Architecture III | 0.000133 | 0.000409 | 0.000979 | 0.001717 |



(a) Cost Plots @ 1 SEC

(b) Cost Plots @ 5 SEC

(c) Cost Plots @ 10 SEC
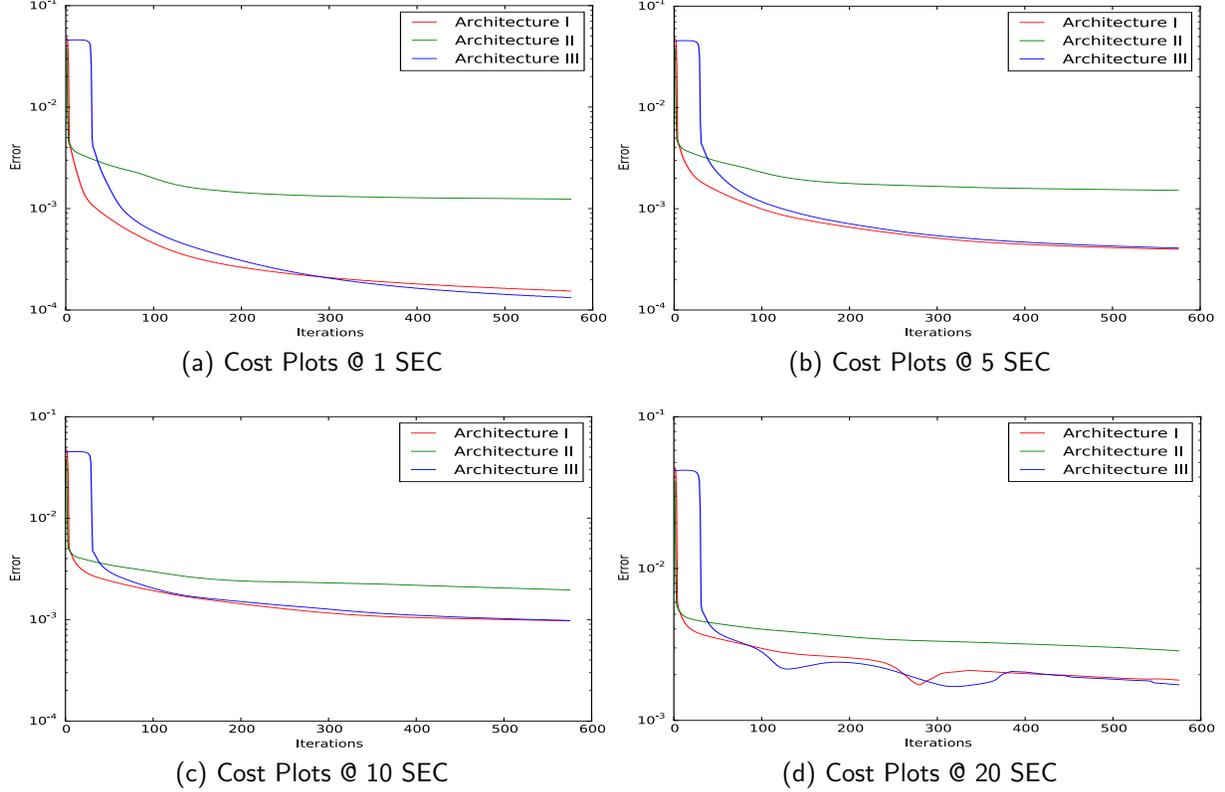
(d) Cost Plots @ 20 SEC

Figure 11: Mean squared error during the training process for the three architectures predicting vibration in 1, 5, 10, and 20 future sec.

## 6.2. Cost Function

Mean squared error was used to train the neural networks as it provides a smoother optimization surface for backpropagation than mean average error. The cost function output for predicting 1 sec, 5 sec, 10 sec and, 20 sec is shown in Figures 11a, 11b, 11c and, 11d respectively. Results are shown in logarithmic scale.

## 6.3. Previous Architecture Results

Mean Squared Error (MSE) (shown in Equation 10) was used as an error measure to train the three architectures, which resulted in values shown in Table 4. Mean Absolute Error (MAE) (shown in Equation 11) was used as a final measure of accuracy for the three architectures, with results shown in Table 5. As the parameters were normalized between 0 and 1, the MAE is also the percentage error.

$$Error = \frac{0.5 \times \sum (Actual\,Vib - Predicted\,Vib)^2}{Testing\,Seconds} \qquad (10)$$

$$Error = \frac{\sum [ABS(Actual\,Vib - Predicted\,Vib)]}{Testing\,Seconds} \qquad (11)$$

Table 4: Previous Testing Process Mean Squared Error

| | Prediction Error | | | |
| | 1 seconds | 5 seconds | 10 seconds | 20 seconds |
| --- | --- | --- | --- | --- |
| Architecture I | 0.000792 | 0.001165 | 0.002926 | 0.010427 |
| Architecture II | 0.010311 | 0.009708 | 0.009056 | 0.012560 |
| Architecture III | 0.000838 | 0.002386 | 0.004780 | 0.041417 |

Table 5: Previous Testing Process Mean Absolute Error

| | Prediction Error | | | |
| | 1 seconds | 5 seconds | 10 seconds | 20 seconds |
| --- | --- | --- | --- | --- |
| Architecture I | 0.028407 | 0.033048 | 0.055124 | 0.101991 |
| Architecture II | 0.098357 | 0.097588 | 0.096054 | 0.112320 |
| Architecture III | 0.027621 | 0.048056 | 0.070360 | 0.202609 |

Figures 12, and Figures 13b, 13d, 13f present the predictions for all the test flights condensed on the same plot. The time shown on the x-axis is the total time for all the test flights. Each flight ends when the vibration reaches the max critical value (normalized to 1) and then the next flight in the test set begins. Figures 14 provides an uncompressed example of Architecture I and Architecture III predicting vibration 5, 10 and, 20 seconds in the future over a single flight from the testing data.

### 6.3.1. Results of Architecture I

The results of this architecture, shown in Table 4, came out to be the best results regarding the overall accuracy of the vibration prediction. There is more misalignment between the actual and calculated vibration values as predictions are made further in the future, as shown in Figures 12a, 12c, 12e, however this is to be expected as it is more challenging to predict further in the future. Also, it can be seen that the prediction of higher peaks is more accurate than the prediction of lower peaks, as if the neural network is tending to learn more about the max critical vibration value, which is favorable for this project.

To test this architecture further, the architecture was trained and tested on the same data set but for predicting vibration just one second in future. As expected, the results showed improvement in mean absolute error over all the test flights by about 0.5% compared to the results of the same architecture predicting five seconds in the future. A plot of the test data prediction for this experiment is shown in Figure 13a. Also, for comparison, a plot for the same flights plotted in Figure 14 is shown in Figure 13g.
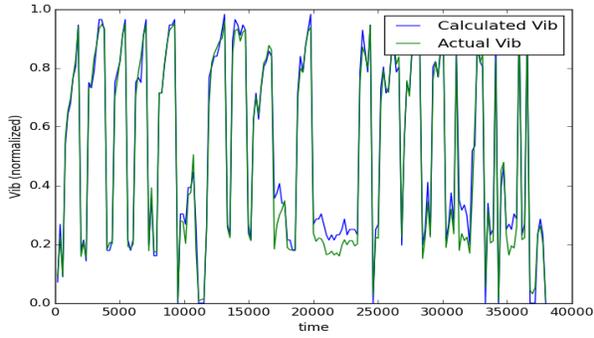
### 6.3.2. Results of Architecture II

The results of this architecture in Table 4 came out to be the least successful in vibration prediction. While it managed to predict much of the vibration, its performance was weak at the peaks (either low or high) compared to the other architectures, as shown in Figures 12d, 12d, 12d. However, it is also worth mentioning that somehow the lower peaks were better at some positions on the curve of this architecture, compared to the other architectures. A potential reason for the poor performance of this architecture is due to using the average of values from the LSTM second layer output, the other two architectures can weight the values from the LSTM second layer output for more accuracy.
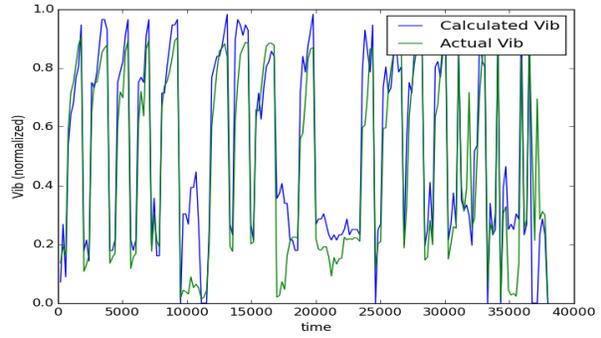
### 6.3.3. Results of Architecture III

This LSTM RNN was one layer deeper and also had 20 seconds memory from the past which was not available for the other two LSTM RNNs used. Although it was the most computationally expensive and had the most chance for deeper learning, the results of this architecture were not as good as expected, as shown in Figure 13. Figures 14b, 14d, 14f provide an uncompressed example of Architecture III predicting vibration 5, 10 and, 20 seconds in the future over a single flight from the testing data. The results of this architecture in Table 4 show that the prediction accuracy for this architecture was less than the more simple Architecture I. While this came counter to the benefits of deeper learning, it does opens door for investigating about how to further tune more complicated LSTM RNNs.
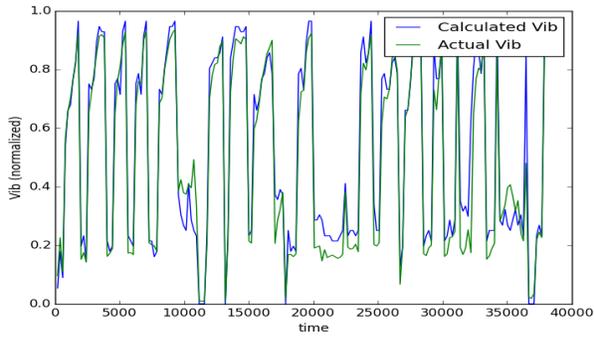
The overall error in Table 4 for the prediction at 20 future seconds was relatively high. Looking at Figure 13f between time 10,000-15,0000, 20,000-25,000 and 35,000-40,000, it can be seen that the calculated
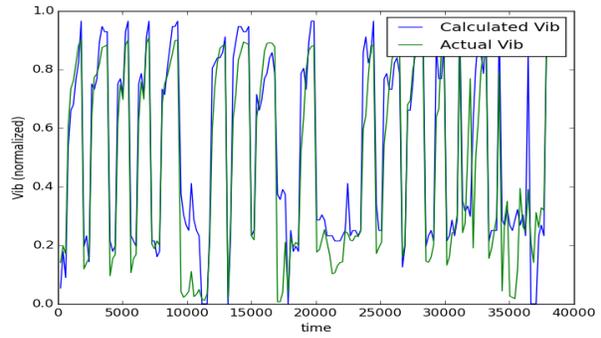
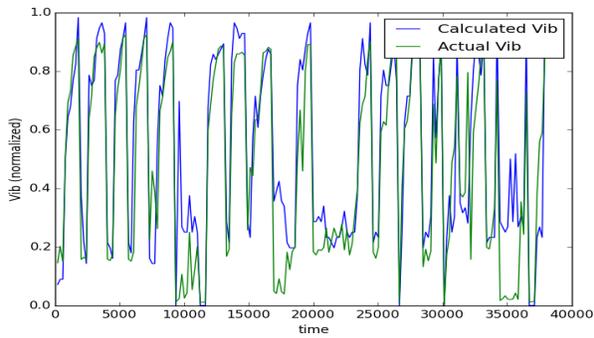(a) Architecture I Results Plot @ 05 SEC
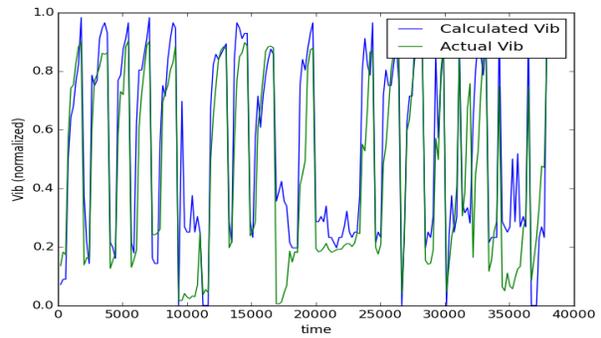
(b) Architecture II Results Plot @ 05 SEC

(c) Architecture I Results Plot @ 10 SEC

(d) Architecture II Results Plot @ 10 SEC

(e) Architecture I Results Plot @ 20 SEC

(f) Architecture II Results Plot @ 20 SEC

Figure 12: Plotted results for Architectures I and II for 5, 10 and 20 seconds in the future.

(a) Architecture I Results Plot @ 1 SEC (all test flight)

(b) Architecture III Results Plot @ 05 SEC (all test flight)

(c) Architecture II Results Plot @ 1 SEC (all test flight)

(d) Architecture III Results Plot @ 10 SEC (all test flight)

(e) Architecture III Results Plot @ 1 SEC (all test flight)
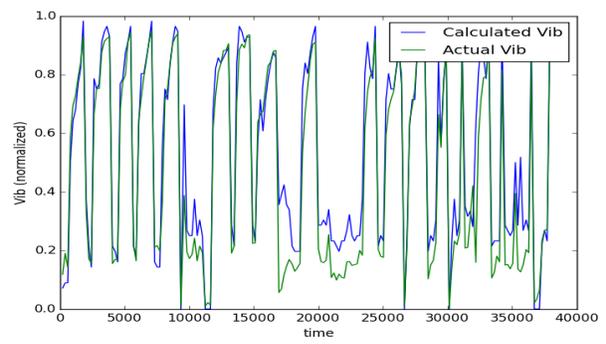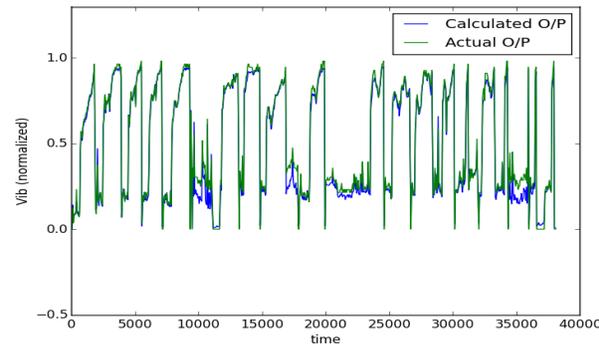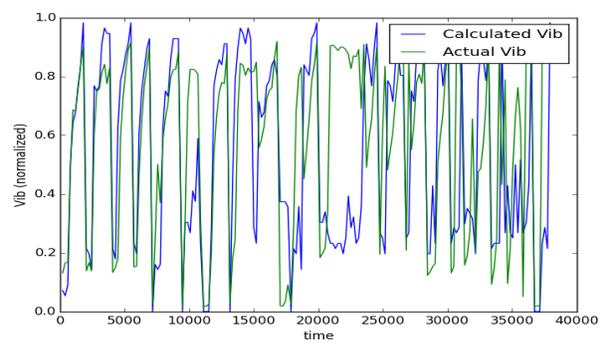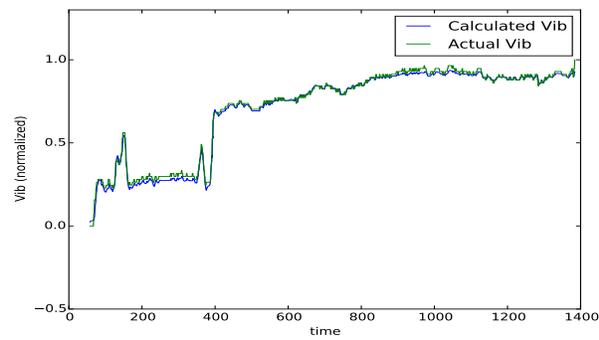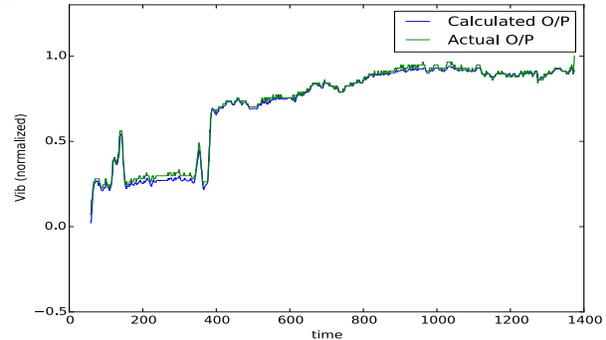
(f) Architecture III Results Plot @ 20 SEC (all test flight)

(g) Architecture I @ 1 SEC (one flight)

(h) Architecture III @ 1 SEC (one flight)

Figure 13: Plotted results (cont.).

curve got very much higher than the actual vibration curve. This strange behavior is unique as it can be seen that the calculated vibration would rarely exceed the actual vibration for all the curves plotted for all the architectures at all scenarios, and it would be for relatively small value if occurred.

This network could potentially gain further improvement if trained for more epochs over the other simpler architectures since it was deeper. This was tried, giving the neural network about double the number of training epochs. However, a significant improvement in the prediction was not achieved. Nonetheless, it is of note that the plots of the cost function of this architecture were not smooth while trained to predict for 20 seconds in the future. This could potentially be a result of under-training or other issues in the training process.

Initially, the training epochs were fixed at 575 for all the architecture as a standard for performance comparison. Further, the performance of this architecture (the mean absolute error) was slightly better than the other architectures when predicting for 1 second in the future. This result supports the believe that this architecture can perform better if given the chance to train for more epochs.

### 6.4. ACO Results

As Architecture I gave the most promising results, it was chosen as the initial candidate for the ACO. The ACO code was run for 1000 iterations with a chosen number of 200 ants. Each run took approximately 4 days. The neural networks were run trained against flights that suffered from the excessive vibration. They were then evaluated against different set of test flights, which also suffered from the same problem. There were 28 flights in the training set, with a total of 41,431 seconds of data. There were 29 flights in the testing set, with a total of 38,126 seconds of data. The networks were allowed to train for 575 epochs to learn and for the cost function output curve to flatten. The minimum value for the pheromones were 1 and the maximum was 20. The pheromones increased only when the network was found to give better fitness than the best fitness in the ACO generated population. The population size was equal to number number of iterations in the ACO process, *i.e.*, the population size was also 1000.
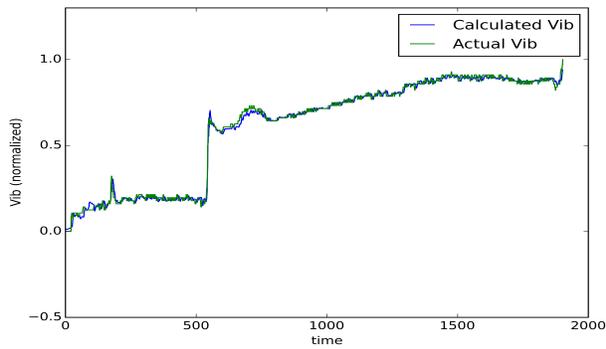
The best version of Architecture I evolved with ACO showed an improvement of 1.35% for predictions 10 seconds in the future, reducing prediction error from 5.51% to 4.17% compared to the architecture's performance before the ACO. Results of the ACO are shown in Figure 15b for Architecture I for a single test flight and they are compared to the performance of the same architecture also for a single test flight. Figures 15c, 15a show the results for the all the test flight for the architecture, before and after the ACO. The plot of the cost function of the training process of the best evolved network is shown in Figure 16.

Returning to an initial question of how the number of the connections in the network affects the soundness of the results, Table 6 shows the top thirty evolved networks with respect to the fitnesses they provide. The table also shows the total number of connections in both $mesh\_1$ (first set of connections in the generated mesh) and $mesh\_2$ (second set of connections in the generated mesh), and the total number weights (connections) in the networks. Comparing these values to the total number of weights in a fully connected Architecture I type network, as shown in Table 2, it is found that total number of weights were reduced by 42% to 45% in the top 30 networks.

The ACO generated mesh (as defined in Algorithm 1) used to generate this topology is shown in the matrices in Equations 12 and 13. It is worth stressing that this topology is not the complete LSTM RNN used in the utilized Architecture I, but rather applies to the individual gates in each cell. Equation 12 is used for any fully connected process and Equation 13 is used for any data-reduction process (This is discussed in details in Section 4).

The topology of the design of the networks' cells are shown in Figures 17 and 18. Figures 17a, 17b, and 17c show the ACO optimized $mesh\_1$, which were used within the "M1" LSTM cells (see Figure 4) in the top three evolved LSTM RNNs. Equation 12 represents $mesh\_1$ of the best evolved neural network which was used to generate Figure 17a.

The evolved networks retained all the elements of $mesh\_2$ , represented by Equation 13, for use in the "M2" LSTM cells (see Figure 5). Figure 18 is used to show this part of the evolved mesh. For clarity, Figure 17b shows the differences between the M1 cells before and after ACO optimization. Figure 17a is simply a LSTM cell "M1" that have its gates' meshes (shown in Figure 17b, Up) substituted with the ACO meshes (shown in Figure 17b, Down). "M2" did not change from its original topology as shown in Figure 5 since all the elements in $mesh\_2$ after the optimization remained ones (Equation 13).

24

(a) ART I predicting vibration 5 seconds in the future for one flight.

(b) ART III predicting vibration 5 seconds in the future for one flight.

(c) ART I predicting vibration 10 seconds in the future for one flight.

(d) ART III predicting vibration 10 seconds in the future for one flight.

(e) ART I predicting vibration 20 seconds in the future for one flight.

(f) ART III predicting vibration 20 seconds in the future for one flight.

Figure 14: Architectures I and III predicting vibration for one flight.

Table 6: ACO Top Thirty Evolved Networks

| No. | Fitness | Number of M1 Connections | Number of M2 Connections | Total Number of Connections | No. | Fitness | Number of M1 Connections | Number of M2 Connections | Total Number of Connections |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.041723 | 139 | 16 | 11,810 | 16 | 0.044603 | 134 | 16 | 11,410 |
| 2 | 0.041927 | 136 | 16 | 11,570 | 17 | 0.044680 | 132 | 16 | 12,050 |
| 3 | 0.042549 | 144 | 16 | 12,210 | 18 | 0.044719 | 140 | 16 | 11,890 |
| 4 | 0.043000 | 137 | 16 | 11,650 | 19 | 0.044838 | 141 | 16 | 11,970 |
| 5 | 0.043134 | 137 | 16 | 11,650 | 20 | 0.044922 | 141 | 16 | 11,970 |
| 6 | 0.043399 | 138 | 16 | 11,730 | 21 | 0.044949 | 143 | 16 | 12,130 |
| 7 | 0.043434 | 137 | 16 | 11,650 | 22 | 0.044982 | 139 | 16 | 11,810 |
| 8 | 0.043456 | 143 | 16 | 12,130 | 23 | 0.045013 | 133 | 16 | 11,330 |
| 9 | 0.043806 | 139 | 16 | 11,810 | 24 | 0.045117 | 137 | 16 | 11,650 |
| 10 | 0.043986 | 144 | 16 | 12,210 | 25 | 0.045240 | 136 | 16 | 11,570 |
| 11 | 0.044109 | 141 | 16 | 11,970 | 26 | 0.045302 | 147 | 16 | 12,450 |
| 12 | 0.044261 | 137 | 16 | 11,650 | 27 | 0.045306 | 144 | 16 | 12,210 |
| 13 | 0.044362 | 143 | 16 | 12,130 | 28 | 0.045363 | 147 | 16 | 12,450 |
| 14 | 0.044483 | 142 | 16 | 12,050 | 29 | 0.045367 | 133 | 16 | 11,330 |
| 15 | 0.044493 | 136 | 16 | 11,570 | 30 | 0.045397 | 146 | 16 | 12,370 |

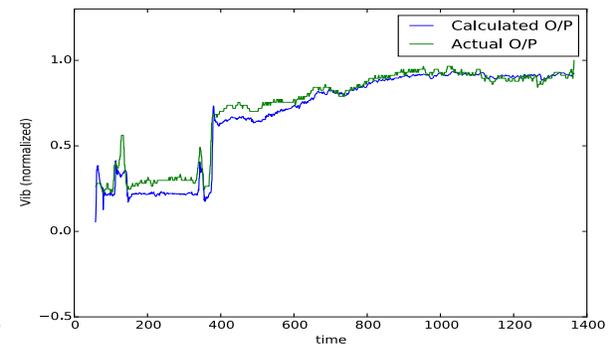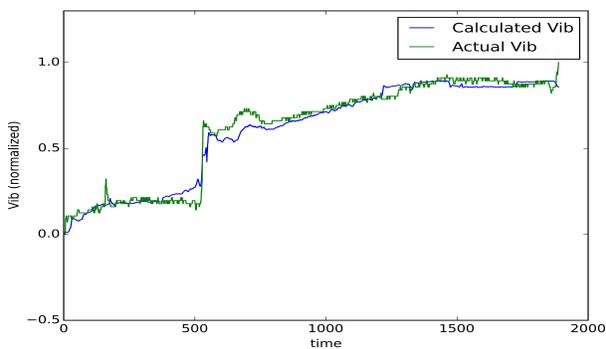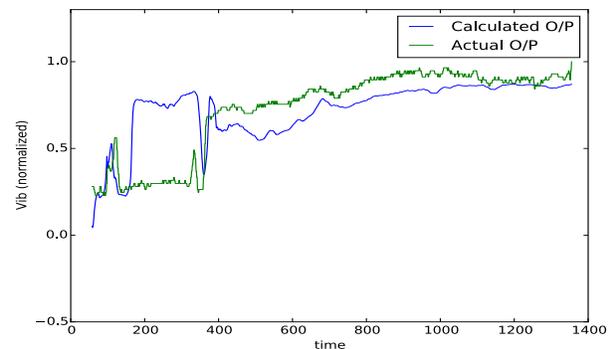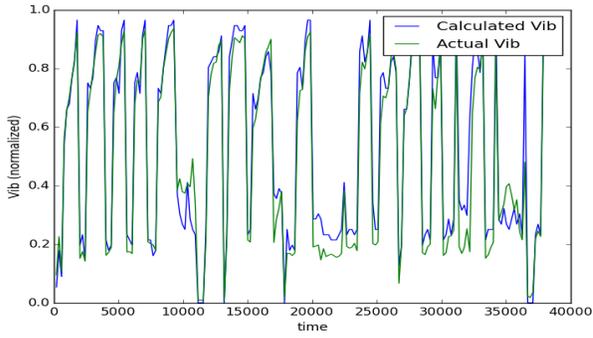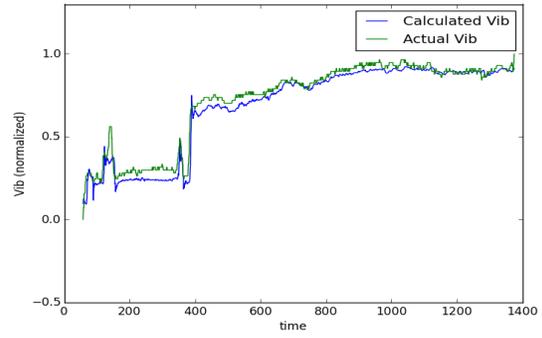The colored nodes in Figure 17a are the input nodes (first line of nodes) at the Main, Input, Forget, and Output gates at the "M1" cells (Figure 4). The diamond nodes in Figure 17a are the hidden layer nodes (second line of nodes) at the Main, Input, Forget, and Output gates at the "M1" cells (Figure 4). The diamond nodes are also the input nodes at the Main, Input, Forget, and Output gates at the "M2" cells (Figure 5). The last single node in Figure 18 is the output of the gates in the "M2" cells.

$$mesh\_1 = \begin{array}{c} i1 \\ i2 \\ i3 \\ i4 \\ i5 \\ i6 \\ i7 \\ i8 \\ i9 \\ i10 \\ i11 \\ i12 \\ i13 \\ i14 \\ i15 \\ bias \end{array} \begin{vmatrix} h1 & h2 & h3 & h4 & h5 & h6 & h7 & h8 & h9 & h10 & h11 & h12 & h13 & h14 & h15 & h16 \\ 0 & 1 & 0 & 1 & 1 & 0 & 1 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 1 & 1 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 1 & 1 & 0 & 1 & 1 & 0 & 1 & 1 & 0 & 1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 & 1 & 0 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 & 1 & 0 & 1 & 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 1 & 0 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 & 1 & 0 & 0 & 0 & 1 & 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 1 & 1 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 & 1 & 1 & 0 & 1 & 0 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 1 \end{vmatrix} \tag{12}$$

$$mesh\_2 = \begin{vmatrix} h1 & h2 & h3 & h4 & h5 & h6 & h7 & h8 & h9 & h10 & h11 & h12 & h13 & h14 & h15 & h16 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{vmatrix} \tag{13}$$

## 7. Discussion and Future Work

The results have shown that the ACO approach for optimizing the gates within LSTM cells can dramatically reduce the number of connections required, while at the same time improve the predictive ability of the recurrent neural network. However, as much of the matrix in Equation 12 is sparse, none of the rows of this matrix had all their elements equal to zero, meaning that none of the inputs ended up being dropped out (which was a potential means for improving predictions, as discussed in Subsection 3.1.2). This is a indication that all the chosen parameters actually had a positive contributing influence on the vibration. On the other hand, it suggests that having additional connections can increase the difficulty of appropriately training the LSTM RNN, resulting in less predictive ability (as in the case of the original unoptimized LSTM RNN architectures).

This approach can open the door to further identify the highest contributors to the vibration problem by examining the number of the connections between the input neurons and the hidden layers, along with the magnitude of those weights. Furthermore, the combined effect of multiple parameters can also be investigated by looking at the input neurons connected to a certain hidden layer's neuron. This would give an idea about the effect those input neurons, which represents the input parameters, have on the vibration as a final output, which could aid in discovering actual cause of the vibration events.

This also opens up the potential for significant future work. While the optimized LSTM RNNs did not drop out any input connections, by increasing the number of flight parameters used as input (even using all available parameters), the algorithm has the potential to determine which parameters contribute most to
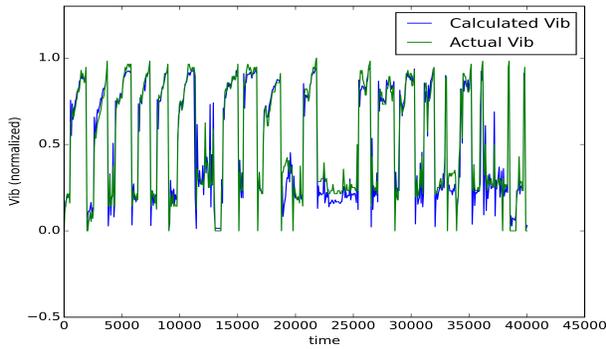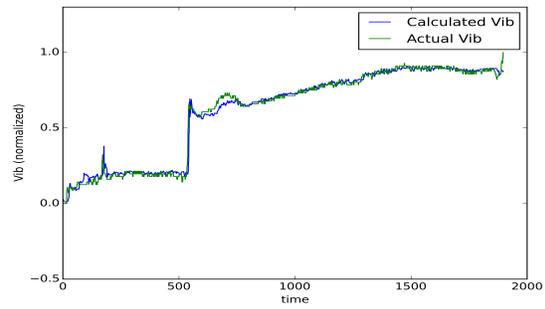
(a) Unoptimized: all test flights

(b) Unoptimized: one test flights

(c) Optimized: all test flights

(d) Optimized: one test flights

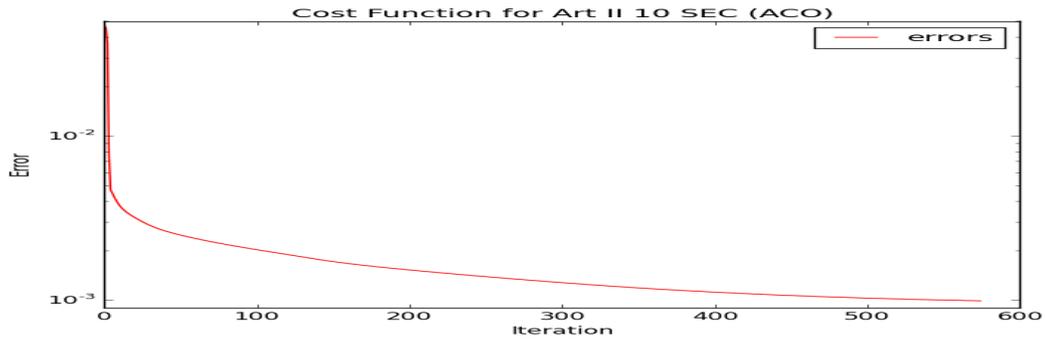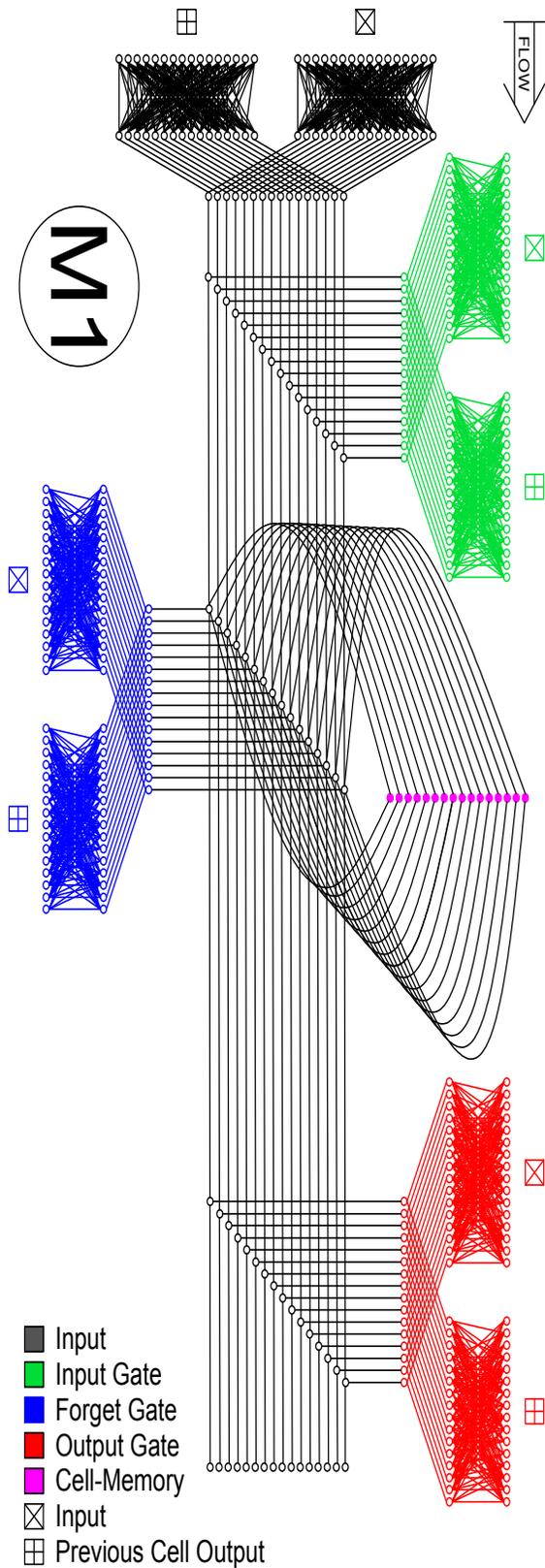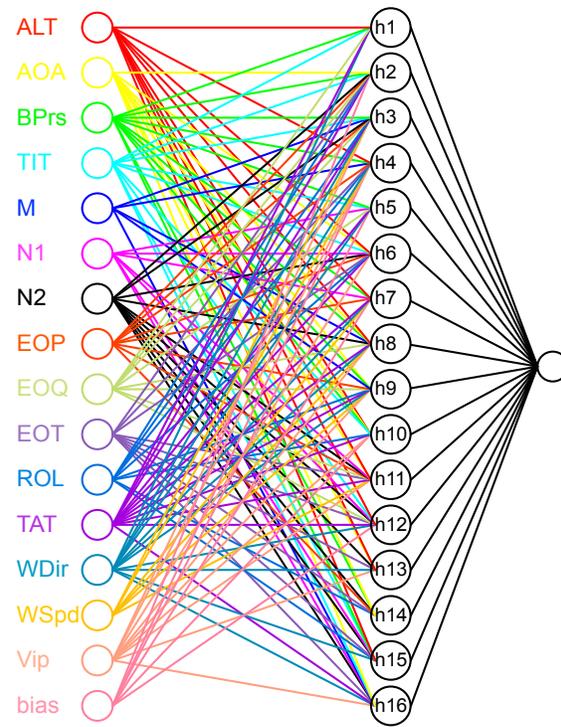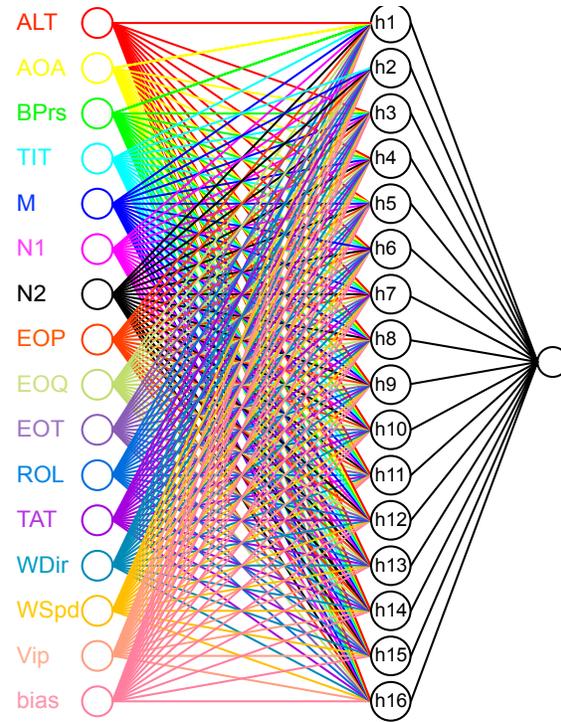Figure 15: Plotted results for predicting ten seconds in the future.



Figure 16: Cost function plot for ACO optimized ART I predicting vibration in 10 future sec

Input
Input Gate
Forget Gate
Output Gate
Cell-Memory
Input
Previous Cell Output

(a) ACO Architecture I First Best Fitness Network: LSRM M1 cell.

(b) Meshes before (up) and after (down) optimization.

(a) ACO Architecture I First Best Fitness Mesh: 155 connections.

(b) ACO Architecture I Second Best Fitness Mesh: 152 connections.

(c) ACO Architecture I Third Best Fitness Mesh: 160 connections.

Figure 17: ACO Architecture I Best Fitness Topologies' Meshes (Equation 12) for at "M1" (Figure 4) LSTM cells: 1000 Iterations, and 200 Ants.

Figure 18: ACO Architecture I Best Fitness Topology's mesh (Equation 13) at "M2" (Figure 5) LSTM cells: 1000 Iterations, and 200 Ants.

the predictive ability, instead of relying on *a priori* expert knowledge to select parameters. Further, in this work one mesh of connections was generated and then used in all the LSTM cell gates at all time-steps. The future work will consider the meshes in the four gates of the LSTM cells at the various LSTM time-steps as variable and will apply the ACO on each of them simultaneously in every ACO iteration.

Future work will also consider optimizing the LSTM RNN structure. The connections of the structure shown in Figures 9 and 7 will be subject to ACO process along with the optimization of the connections within the LSTM cells. This has the potential to make a large step forward in the evolution of the LSTM RNNs as it will allow for connections between non-adjacent cells, and potentailly even dropping out certain unused cells and potentailly even full layers from the LSTM RNN.

Lastly, work investigating the tuning of the ACO hyperparameters can be done to improve how quickly the algorithm converges to optimal LSTM RNN structures. For example, modifying the number of ants, reducing pheromones on paths from LSTM RNNs with lower fitness, and periodically refreshing the pheromones levels by decreasing all of its levels by certain amount.

## Acknowledgments

## References

## References

[1] A. V. Srinivasan, *The American Society of Mechanical Engineers* **1997**,

[2] Hochreiter, S.; Schmidhuber, J. *Neural computation* **1997**, *9*, 1735–1780.

[3] Di Persio, L.; Honchar, O.

[4] S. Hochrieter & J. Schmidhuber, *Neural Computation 9(8):1735-1780* **1997**,

[5] Felder, M.; Kaifel, A.; Graves, A. Wind power prediction using mixture density recurrent neural networks. Poster Presentation gehalten auf der European Wind Energy Conference. 2010.

[6] Choi, E.; Bahadori, M. T.; Sun, J. *arXiv preprint arXiv:1511.05942* **2015**,

[7] Maknickienė, N.; Maknickas, A. Application of neural network for forecasting of exchange rates and forex trading. The 7th international scientific conference" Business and Management. 2012; pp 10–11.

[8] Yao, X. *Proceedings of the IEEE* **1999**, *87*, 1423–1447.

[9] Siebel, N. T.; Botel, J.; Sommer, G. Efficient neural network pruning during neuro-evolution. Neural Networks, 2009. IJCNN 2009. International Joint Conference on. 2009; pp 2920–2927.

[10] Schmidhuber, J. *Neural networks* **2015**, *61*, 85–117.

[11] Zhang, B.-T.; Muhlenbein, H. *Complex systems* **1993**, *7*, 199–220.

[12] Kohl, N. F. **2009**,

[13] Whiteson, S. Improving reinforcement learning function approximators via neuroevolution. Proceedings of the fourth international joint conference on Autonomous agents and multiagent systems. 2005; pp 1386–1386.

[14] Floreano, D.; Dürr, P.; Mattiussi, C. *Evolutionary Intelligence* **2008**, *1*, 47–62.

[15] Turner, A. J.; Miller, J. F. *Research and Development in Intelligent Systems XXX*; Springer, 2013; pp 213–226.

[16] Desell, T.; Clachar, S.; Higgins, J.; Wild, B. Evolving Deep Recurrent Neural Networks Using Ant Colony Optimization. European Conference on Evolutionary Computation in Combinatorial Optimization. 2015; pp 86–98.

[17] ElSaid, A.; Wild, B.; Higgins, J.; Desell, T. Using LSTM recurrent neural networks to predict excess vibration events in aircraft engines. e-Science (e-Science), 2016 IEEE 12th International Conference on. 2016; pp 260–269.

[18] Dorigo, M.; Maniezzo, V.; Colorni, A. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* **1996**, *26*, 29–41.

[19] Dorigo, M.; Gambardella, L. M. *IEEE Transactions on evolutionary computation* **1997**, *1*, 53–66.

[20] Bianchi, L.; Gambardella, L. M.; Dorigo, M. An ant colony optimization approach to the probabilistic traveling salesman problem. International Conference on Parallel Problem Solving from Nature. 2002; pp 883–892.

[21] Manfrin, M.; Birattari, M.; Stützle, T.; Dorigo, M. Parallel ant colony optimization for the traveling salesman problem. International Workshop on Ant Colony Optimization and Swarm Intelligence. 2006; pp 224–234.

[22] K. Socha, *Future Generation Computer Systems* **2004**, 25–36.

[23] K. Sochaand & M. Dorigo, *Journal of Petroleum Science and Engineering, 77(3):375385* **2011**,

[24] Socha, K. *Ant colony optimisation for continuous and mixed-variable domains*; VDM Publishing Saarbrücken, 2009.

[25] Bilchev, G.; Parmee, I. *Evolutionary Computing* **1995**, 25–39.

[26] N. Monmarch and G. Venturini and M. Slimane, *Future Generation Computer Systems* **2000**, *16*, 937–946.

[27] Dréo, J.; Siarry, P. A new ant colony algorithm using the heterarchical concept aimed at optimization of multiminima continuous functions. International Workshop on Ant Algorithms. 2002; pp 216–221.

[28] R. Ashena and J. Moghadasi, *European journal of operational research, 185(3):11551173* **2008**,

[29] Blum, C.; Socha, K. Training feed-forward neural networks with ant colony optimization: An application to pattern classification. Fifth International Conference on Hybrid Intelligent Systems (HIS'05). 2005; pp 6–pp.

[30] J.-B. Li and Y.-K. Chung, *In Transmission and Distribution Conference and Exhibition: Asia and Pacific, 2005 IEEE/PES, pages 15. IEEE* **2005**,

[31] M. Unal, M. Onat, and A. Bal, *In Signal Processing and Communications Applications Conference (SIU), 2010 IEEE 18th, pages 471474. IEEE* **2010**,

[32] Moubray, J. *Reliability-centered maintenance*; Industrial Press Inc., 1997.

[33] Chatfield, C. *The analysis of time series: an introduction*; CRC press, 2016.

[34] Boukary, N. A. A COMPARISON OF TIME SERIES FORECASTING LEARNING ALGORITHMS ON THE TASK OF PREDICTING EVENT TIMING. Ph.D. thesis, Royal Military College of Canada, 2016.

[35] Goel, H.; Melnyk, I.; Oza, N.; Matthews, B.; Banerjee, A.

[36] Nairac, A.; Townsend, N.; Carr, R.; King, S.; Cowley, P.; Tarassenko, L. *Integrated Computer-Aided Engineering* **1999**, *6*, 53–66.

[37] Clifton, D. A.; Bannister, P. R.; Tarassenko, L. A framework for novelty detection in jet engine vibration data. Key engineering materials. 2007; pp 305–310.

[38] Gers, F. A.; Schraudolph, N. N.; Schmidhuber, J. *Journal of machine learning research* **2002**, *3*, 115–143.

[39] Desell, T.; Clachar, S.; Higgins, J.; Wild, B. In *Parallel Problem Solving from Nature - PPSN XIII*; Bartz-Beielstein, T., Branke, J., Filipič, B., Smith, J., Eds.; Lecture Notes in Computer Science; Springer International Publishing, 2014; Vol. 8672; pp 771–781.

[40] Kennedy, J. *Encyclopedia of machine learning*; Springer, 2011; pp 760–766.

[41] Poli, R.; Kennedy, J.; Blackwell, T. *Swarm intelligence* **2007**, *1*, 33–57.

[42] Storn, R.; Price, K. *Journal of global optimization* **1997**, *11*, 341–359.

[43] Blum, C.; Li, X. *Swarm Intelligence*; Springer, 2008; pp 43–85.

[44] Dorigo, M.; Birattari, M.; Stutzle, T. *IEEE computational intelligence magazine* **2006**, *1*, 28–39.

[45] Dorigo, M.; Stützle, T. *Handbook of metaheuristics*; Springer, 2010; pp 227–263.

[46] M. Dorigo and L. M. Gambardella, *BioSystems, 43(2):7381* **1997**,

[47] Felix A. Gers, Jrgen Schmidhuber, and Fred Cummins, *Neural Computation, Vol. 12, No. 10 , Pages 2451-2471* **2000**,

[48] Gers, F. A.; Eck, D.; Schmidhuber, J. *Neural Nets WIRN Vietri-01*; Springer, 2002; pp 193–200.

[49] Donahue, J.; Anne Hendricks, L.; Guadarrama, S.; Rohrbach, M.; Venugopalan, S.; Saenko, K.; Darrell, T. Long-term recurrent convolutional networks for visual recognition and description. Proceedings of the IEEE conference on computer vision and pattern recognition. 2015; pp 2625–2634.

[50] Chao, L.; Tao, J.; Yang, M.; Li, Y.; Wen, Z. *arXiv preprint arXiv:1603.08321* **2016**,

[51] Eck, D.; Schmidhuber, J. *Istituto Dalle Molle Di Studi Sull Intelligenza Artificiale* **2002**, *103*.

[52] Stanley, K. O.; Miikkulainen, R. *Evolutionary computation* **2002**, *10*, 99–127.

[53] Annunziato, M.; Lucchetti, M.; Pizzuti, S. *Proc. EUNITE02, Albufeira, Portugal* **2002**,

[54] Larochelle, H.; Bengio, Y.; Louradour, J.; Lamblin, P. *Journal of Machine Learning Research* **2009**, *10*, 1–40.

[55] Kandel, E. R.; Schwartz, J. H.; Jessell, T. M.; Siegelbaum, S. A.; Hudspeth, A. J. *Principles of neural science*; McGraw-hill New York, 2000; Vol. 4.

[56] Lewis, J. P. Fast normalized cross-correlation. Vision interface. 1995; pp 120–123.

[57] Dalcín, L.; Paz, R.; Storti, M.; DElía, J. *Journal of Parallel and Distributed Computing* **2008**, *68*, 655–662.

[58] Szymanski, B.; Desell, T.; Varela, C. The Effect of Heterogeneity on Asynchronous Panmictic Genetic Search. Proc. of the Seventh International Conference on Parallel Processing and Applied Mathematics (PPAM'2007). Gdansk, Poland, 2007.

[59] Desell, T.; Anderson, D.; Magdon-Ismail, M.; Heidi Newberg, B. S.; Varela, C. An Analysis of Massively Distributed Evolutionary Algorithms. The 2010 IEEE congress on evolutionary computation (IEEE CEC 2010). Barcelona, Spain, 2010.

[60] Theano Development Team, *arXiv e-prints* **2016**, *abs/1605.02688*.