

Railway Track Circuit Fault Diagnosis Using Recurrent Neural Networks

Tim de Bruin, Kim Verbert, and Robert Babuška

Abstract—Timely detection and identification of faults in railway track circuits are crucial for the safety and availability of railway networks. In this paper, the use of the long-short-term memory (LSTM) recurrent neural network is proposed to accomplish these tasks based on the commonly available measurement signals. By considering the signals from multiple track circuits in a geographic area, faults are diagnosed from their spatial and temporal dependences. A generative model is used to show that the LSTM network can learn these dependences directly from the data. The network correctly classifies 99.7% of the test input sequences, with no false positive fault detections. In addition, the t-Distributed Stochastic Neighbor Embedding (t-SNE) method is used to examine the resulting network, further showing that it has learned the relevant dependences in the data. Finally, we compare our LSTM network with a convolutional network trained on the same task. From this comparison, we conclude that the LSTM network architecture is better suited for the railway track circuit fault detection and identification tasks than the convolutional network.

Index Terms—Fault diagnosis, long-short-term memory (LSTM), recurrent neural network (RNN), track circuit.

I. INTRODUCTION

AS RAILWAY networks are becoming busier, they are required to operate with increasing levels of availability and reliability [1]. To enable the safe operation of a railway network, it is crucial to detect the presence of trains in the sections of a railway track. The railway track circuit is worldwide the most commonly used component for train detection. To prevent accidents, the detection system is designed to be fail safe, meaning that in the case of a fault, the railway section is reported as occupied.

When this happens, trains are no longer allowed to enter the particular section. This avoids collisions, but leads to train delays. Moreover, in spite of the fail-safe design of the track circuit, there are situations in which the railway section can be incorrectly reported as free, which can potentially lead to dangerous situations. Therefore, to guarantee both safety and a high availability of the railway network, it is very important to prevent track circuit failures. This requires a

preventive maintenance strategy to ensure that the components are repaired or replaced before a fault develops into a failure. To schedule the maintenance of the track circuits in the most efficient and effective manner, it is necessary to detect and identify the faults as soon as possible.

In this paper, we propose a neural network approach to fault diagnosis in railway track circuits. The fault diagnosis task comprises the detection of faulty behavior and the determination of the cause of that behavior.

Since the railway track circuit network is a large network, it is not realistic to assume that additional monitoring devices will be installed on each track circuit. Therefore, this paper assumes only the availability of data that are currently measured in track circuits. By analyzing the measurement signals from several track circuits in a small area over time, the fault cause can be inferred from the spatial and temporal dependences [2]. In contrast to [2], in this paper, a data-based approach to fault diagnosis is considered, namely, an artificial recurrent neural network (RNN) called the long-short-term memory (LSTM) network [3].

Artificial neural networks have recently achieved state-of-the-art performance on a range of challenging pattern recognition tasks, such as image classification [4] and speech recognition [5]. Some of the advances made in these domains can be applied to fault diagnosis problems as well, which makes the use of neural networks an interesting option in this domain.

Learning the long-term temporal dependences that are characteristic of the faults in the track circuit case presents a challenge to standard neural networks. The LSTM network deals with this problem by introducing memory cells into the network architecture.

Currently, not enough measurement data are available to train the network and to verify its performance. Therefore, we have combined the available data with qualitative knowledge of the fault behaviors [2], and we have constructed a generative model. The performance of the proposed approach is demonstrated using synthetic data produced by this model. However, as the amount of available track circuit data is expected to increase rapidly over time, we expect that the method will be relevant.

A. Related Work

Several methods for fault diagnosis in railway track circuits have been proposed in the literature [1], [2], [6]–[10]. A distinction can be made between methods that use the data collected by a measurement train [6], [7], [9], [10]

Manuscript received March 31, 2015; accepted April 1, 2016. This research is part of the STW/ProRail project “Advanced monitoring of intelligent rail infrastructure (ADMIRE)”, project 12235, supported by the Dutch Technology Foundation STW. It is also part of the research programme Deep Learning for Robust Robot Control (DL-Force) with project number 656.000.003. Both projects are partly financed by the Netherlands Organisation for Scientific Research (NWO). (Corresponding author: Tim de Bruin.)

The authors are with the Delft Center for Systems and Control, Delft University of Technology, Delft 2628 CD, The Netherlands (e-mail: t.d.debruin@tudelft.nl; k.a.j.verbert@tudelft.nl; r.babuska@tudelft.nl).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TNNLS.2016.2551940

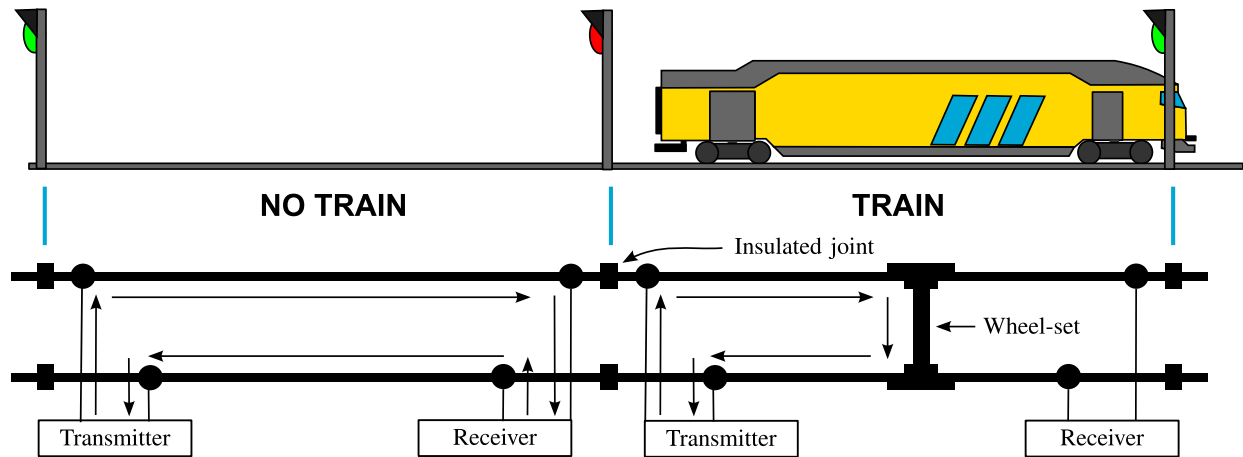


Fig. 1. Current flow in a track circuit. Each track circuit detects the absence of trains in a section of a railway track. Subsequent sections are separated from each other by insulated joints.

and methods that use data collected by track-side monitoring devices [1], [2], [8]. In this paper, the track-side monitoring devices are considered, because they continuously monitor the system health and are, therefore, suitable for the early diagnosis of faults. The main difference compared with the approaches in [1] and [8] is that in those works, multiple monitoring signals are used, while in this paper, for each track circuit, only one measurement signal is available. The main difference compared with the approach in [2] is that in [2], a knowledge-based approach is proposed, while we consider a data-based approach, namely, an LSTM network.

The use of spatial fault dependences for the diagnosis of faults is relatively new to the railway track circuit setting [2], although it is more commonly used in other domains [11]–[13].

To the best of our knowledge, LSTM networks have not been previously proposed for fault diagnosis in railway track circuits. However, many applications of neural networks to fault diagnosis and condition monitoring problems can be found in the literature. One popular approach is to use a deep belief network [14]. The stochastic nature of these networks make them a natural fit to fault detection. By training exclusively on the examples from healthy behavior, the network can determine the probability that a new input vector does not come from the class of healthy states.

One example of this principle is given in [15], where a deep belief network is trained to detect faults in electric motors. In [16], a deep belief network is used to create an industrial soft sensor. The network predicts the value of a process variable based on the values of many other variables. However, it does not take the temporal developments of these variables into account. When these methods do take a time sequence as an input, they often consider a sequence of fixed length. In contrast, we use a recurrent network, which allows the predictions of the network to be updated at every input time step while keeping a memory of the past inputs.

Methods using RNNs have also been discussed in the literature. An example closely related to this paper is given in [17], where echo state networks are trained to learn the

spatial and temporal dependences in a distributed sensor network. Faults are detected by predicting the values that the sensors will measure and comparing these with the true values. Methods for fault classification based on predicting the output of a system are common as well. One example is [18], in which for each fault category, a separate RNN model predicts the output of the system given the inputs. The fault is then identified by determining which model best explains the measured outputs. In contrast to these methods, our method learns to detect and classify faults directly from the measurements. In addition, using the LSTM network architecture allows us to learn longer term temporal dependences.

The rest of this paper is organized as follows. In Section II, the working of a track circuit is discussed. In Section III, the structure and the working of the LSTM Network that is used to identify the faults are discussed. The results of using the proposed neural network with the synthetic data are given in Section IV, together with an analysis of the trained network using the visualization method t-SNE [19]. In Section V, a comparison is made between the proposed LSTM network and a convolutional network. The conclusions of this paper are given in Section VI. In Appendix A, a number of faults that can cause a track circuit to fail are presented, with special attention given to the spatio-temporal dependences that make it possible to identify these faults from the measured or generated data. Appendix B describes the generative model that is used to produce the training and test data.

II. TRACK CIRCUITS

To enable the safe operation of a railway network, track circuits are used to detect the absence of a train in a section of railway track. Trains are only allowed to enter track sections, which the corresponding track circuit has reported to be free.

A track circuit works by using the rails in a track section as conductors that connect a transmitter at one end of the section to a receiver at the other end, as shown in Fig. 1. When no train is present in the section, the transmitter will energize a relay in the receiver, which indicates that the section is free.

When a train enters the section, the wheel sets of the train form a short circuit, as shown in Fig. 1. This causes the current flow through the receiver to decrease to a level, where the relay is no longer energized and the section is reported as occupied.

The correct operation of a track circuit depends on the electrical current through the receiver. In the absence of a train in the section, the current must be high enough to energize the relay. Conversely, in the presence of a train, the current must be low enough, so that the relay is de-energized. To maintain the safety and availability of the railway network, it is important to detect all possible faults in the system. Moreover, to schedule preventive maintenance on the track circuits, it is important to identify the fault type and to determine the development of the fault severity over time.

A. Fault Diagnosis

Every track circuit has different electrical properties which results in different values of the high current $I_h(t)$ when no train is present, and of the low current $I_l(t)$ when a train is present. In addition, the transients between these values may be different. The current levels also depend on environmental influences and on the properties of the train passing through the section. For these reasons, it is not possible to adequately detect the presence of a fault by only considering the electrical current $I(t)$ during the passing of a single train. In this paper, we consider the current signals from several track circuits in the same geographic area, measured over a longer period of time. This makes it possible to not only detect the presence of a fault, but to also distinguish between different fault types. The reasoning behind this approach is that different faults have different spatial and temporal footprints [2]. The faults that are considered in this paper are as follows:

- 1) insulated joint defect;
- 2) conductive object (across the insulated joints);
- 3) mechanical rail defect;
- 4) electrical disturbance;
- 5) ballast degradation.

A description of these fault types, together with their spatial and temporal footprints, is given in Appendix A.

B. Generative Model

To enable the development, testing, and comparison of the condition monitoring methods, we have developed a generative model. This model is based on a qualitative understanding of the system and the effect of the faults considered, as well as on limited set of measurement data available from real-world track circuits. This model, together with a strategy for sampling the electrical current, is described in Appendix B.

III. NEURAL NETWORK

Artificial neural networks have achieved the state-of-the-art performance on several pattern recognition tasks. One reason for these successes is the use of a strategy called end-to-end learning. This strategy is based on moving away from hand-crafted feature detectors and manually integrating prior knowledge into the network. Instead, networks are trained

to produce their end results directly from the raw input data. To use end-to-end learning, a large labeled data set is required. When this requirement is met, the benefits of a holistic learning approach tend to be larger than the benefits of explicitly using prior knowledge [20].

One example of a field in which this strategy has been successfully applied is image recognition. On this problem, convolutional networks achieve the state-of-the-art performance by using raw pixel values, instead of using hand-crafted feature detectors as inputs [4]. Another example is speech recognition, in which methods using phonemes as an intermediate representation are being replaced by the methods transcribing sound data directly into letters [5].

For the track circuit fault diagnosis case, there are currently not enough labeled data available. However, the measuring equipment that records these data has been installed. Therefore, it is reasonable to assume that at some future time, the data requirement will be met. The neural network proposed in this paper is trained and tested with synthetic data from our generative model. This enables us to analyze the opportunities of applying end-to-end learning to the track circuit fault diagnosis problem.

A. Network Architecture

The prior knowledge of the spatial and temporal fault dependences will not be explicitly integrated into the neural network. It is, however, important to give the network a structure that enables it to learn these dependences from the data.

In order to take the spatial dependences into account, the network input consists of the electrical current signals from five separate track circuits. The signals come from the track circuit that is being diagnosed $I_B(t)$, as well as two other track circuits on the same track $\{I_A(t), I_C(t)\}$ and two track circuits on an adjacent track $\{I_D(t), I_E(t)\}$.

For detecting temporal dependences, an RNN is a natural choice, since the recurrent connections in the network allow it to store memories of past events. However, standard RNNs struggle to learn the long-term time dependences. This is due to the vanishing gradient problem [3]. A popular solution to this problem is the use of the LSTM network architecture.

1) *LSTM Cell*: LSTM networks are able to learn long-term time dependences by introducing specialized memory cells into the network architecture. The structure of the memory cell is shown in Fig. 2. The units a and b are the input and output units, respectively. M is the memory unit. It can remember a value through a recurrent connection with itself. The neurons denoted by g are the gate units. The input gate i determines when a new input is added to the value of the memory unit by multiplying the output of the input unit a by the output of the gate unit. In a similar way, the forget gate f determines when the value in the memory unit is kept constant and when it is reduced or reset. The output gate o determines when the cell outputs its value.

Our network has two hidden layers containing 250 LSTM cells each. This configuration was empirically found to reliably

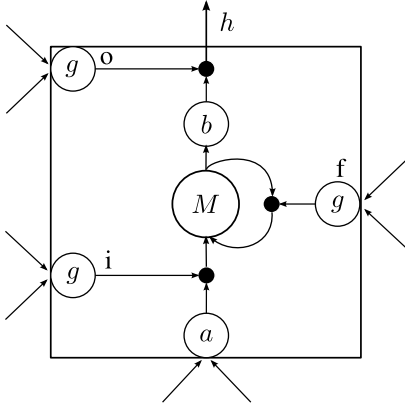


Fig. 2. Architecture of the LSTM memory cell. Black dots: multiplication of the outputs of the gate units g by the outputs of the regular units.

yield good results for this problem. Smaller networks resulted in worse performance and larger networks did not improve the performance further, while requiring significantly increased training times. In general, the ideal size of the network is based on the complexity of the problem, the amount of available training data, and the available computational resources.

The inputs to each LSTM cell j in layer l consist of the inputs to the layer at that time step $x^l(T)$, as well as the outputs of all LSTM cells in layer l at the previous time step $h^l(T-1)$. The equations that describe LSTM cell j in layer l are

$$i_j^l(T) = \text{sigm}(W_{xi_j^l} x^l(T) + W_{hi_j^l} h^l(T-1) + b_{i_j^l}) \quad (1)$$

$$f_j^l(T) = \text{sigm}(W_{xf_j^l} x^l(T) + W_{hf_j^l} h^l(T-1) + b_{f_j^l}) \quad (2)$$

$$a_j^l(T) = \text{tanh}(W_{xa_j^l} x^l(T) + W_{ha_j^l} h^l(T-1) + b_{a_j^l}) \quad (3)$$

$$o_j^l(T) = \text{sigm}(W_{xo_j^l} x^l(T) + W_{ho_j^l} h^l(T-1) + b_{o_j^l}) \quad (4)$$

$$M_j^l(T) = f_j^l(T)(M_j^l(T-1)) + i_j^l(T)a_j^l(T) \quad (5)$$

$$h_j^l(T) = o_j^l(T)\text{tanh}(M_j^l(T)). \quad (6)$$

2) *Inputs and Outputs*: For each of the five track circuits in Fig. 3, the current magnitude is sampled four times during a train-passing event. The details of this sampling procedure are described in Appendix B. The resulting 20 current values for each train-passing event T are the inputs to the first hidden layer for that train-passing event time step: $x^1(T) = [I_A^1(T) \dots I_E^4(T)]$.

The outputs of the first hidden layer are the inputs of the second hidden layer: $x^2(T) = h^1(T)$. The outputs of the second hidden layer are the inputs to the output layer of the network. This layer consists of six softmax classification units; one for the healthy state and five for each of the fault categories. They give the likelihood that the network assigns to each category c at time step T as

$$P(Y=c)(T) = \frac{e^{W_c h^2(T) + b_c}}{\sum_{d=1}^6 e^{W_d h^2(T) + b_d}}. \quad (7)$$

A complete overview of the network is shown in Fig. 3.

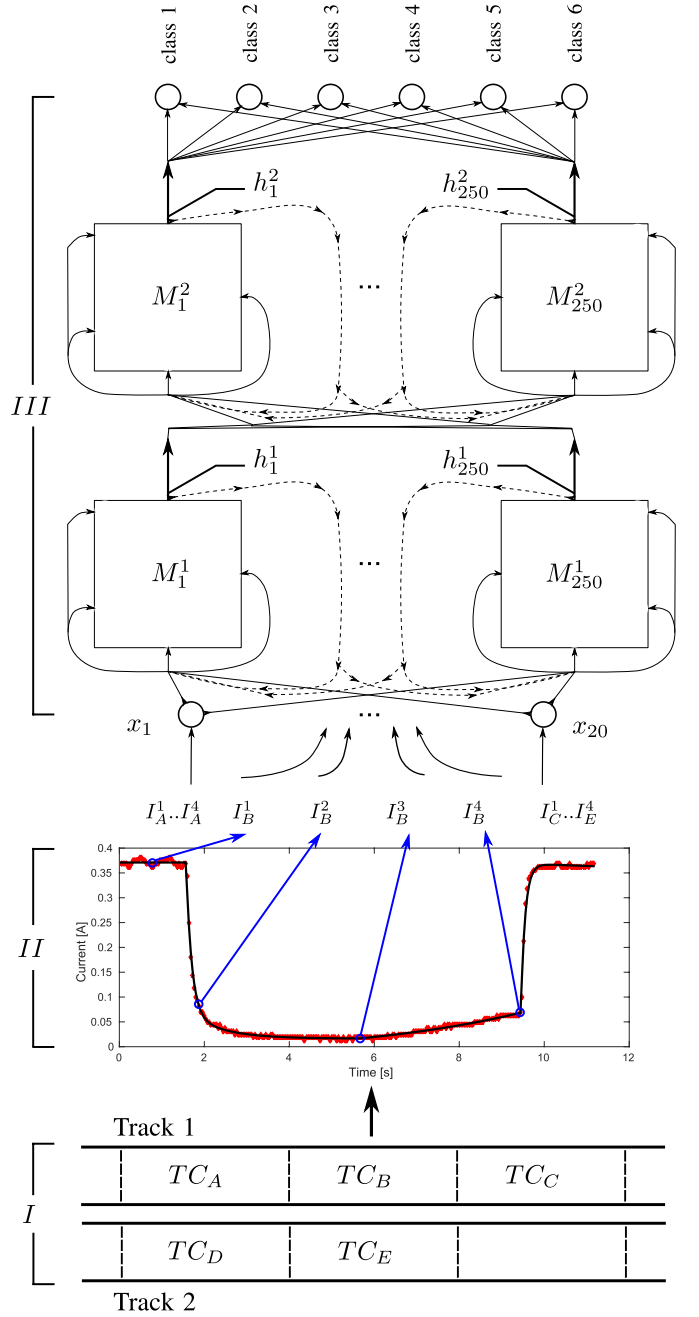


Fig. 3. Fault diagnosis process overview. For each train-passing event T , the current time sequence of the five track circuits (I) is sampled (II). These samples are the input to the neural network (III), which uses them to update the likelihood $P(Y=c)(T)$ of the six different fault classes.

B. Network Training

To train the neural network, two data sets are generated. The first one is a training data set with 21 600 sequences. The second is a validation data set containing 600 sequences. For each sequence, the properties of the track circuits and the properties of the fault are stochastically determined. Each sequence has a length of 2000 train-passing events. This relates to a time period of 100 days. Note that although more trains are likely to pass through the considered sections, it is important to keep the temporal dependences from becoming too long term. Therefore, it might be necessary to limit the number

of train-passing events per day that are used as network inputs.

The network is trained to give a classification of the sequence at every time step \mathcal{T} . The target for this classification $t(\mathcal{T})$ is the healthy state, unless the sequence contains a fault for which the severity at that time step \mathcal{T} is above 0.15. The severity of the fault is between 0 and 1. A fault severity of 0 will have no influence on the electrical current levels, and a fault severity of 1 will influence the current enough to cause a failure, where the track circuit is no longer able to function correctly. The value of 0.15 is chosen to detect the faults as early as possible without having any false positive fault detections. Based on the target classifications $t(\mathcal{T})$, the network is trained to minimize the negative log likelihood loss function

$$l(\mathcal{T}) = -\log(P(Y(\mathcal{T}) = t(\mathcal{T}))). \quad (8)$$

The network is trained with the backpropagation through time algorithm [21] on the sequences in the training data set. The network is unrolled for 500 time steps. First, the network activations and outputs are calculated for these 500 time steps. Then, moving backward through time, the error gradients are calculated and the weights are updated. Finally, the activations of the network at the final time step are used as the initial network activations for the subsequent subsequence of 500 time steps. This process is repeated until all 2000 time steps in the sequence are processed. To improve the efficiency, 56 sequences are processed simultaneously in a minibatch using stochastic gradient descent.

During the training on the training data set, the performance according to (8) on the validation data set is monitored. When this performance stops improving, the learning rate is lowered. After the training is complete, the network weights that resulted in the best performance on the validation data set are used to test the network.

IV. RESULTS

To test the trained network, a test data set is generated containing 1500 sequences.

A. Prediction Accuracy

To test the performance of the network, the test data set is presented to the network. At the final time step of the sequences, the class that is assigned the highest probability is compared with the correct diagnosis for that time step.

Of the 1500 sequences, 1495 were identified correctly. The confusion matrix is given in Table I. An example of a complete input sequence with the resulting classification outputs is shown in Fig. 4, from which it can be seen that the network is insensitive to current drops that are not caused by faults and assigns the majority of the probability to the correct category exactly according to the trained target classifications for each time step. This shows that faults can not only be classified correctly but also identified in a timely fashion, long before they lead to a failure.

TABLE I

CONFUSION MATRIX FOR THE FAULT DIAGNOSIS TASK ON THE TEST DATA SET WITH 1500 SEQUENCES. THE ROWS INDICATE THE TRUE CLASS AND THE COLUMNS REPRESENT THE PREDICTED CLASS

true cat. / pred cat.	1	2	3	4	5	6
1 (healthy)	754	0	0	0	0	0
2 (Insulated joint defect)	0	131	0	1	0	0
3 (conductive object)	1	0	238	0	0	0
4 (mechanical rail defect)	0	1	0	249	0	0
5 (electrical disturbance)	0	0	0	0	4	0
6 (Ballast degradation)	2	0	0	0	0	119

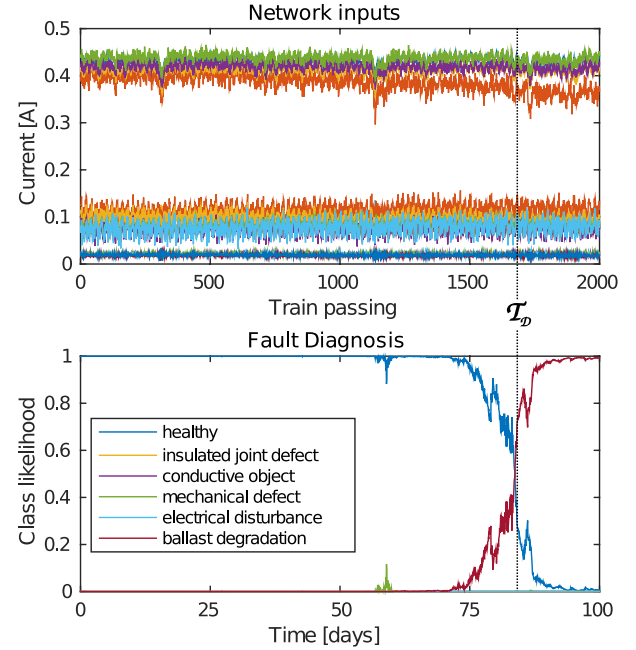


Fig. 4. Network inputs and output for one realization of a ballast degradation fault sequence. The detection time \mathcal{T}_D marks the detection threshold. Before this point, the correct classification is healthy, and after this point, the correct classification is ballast degradation.

B. Misclassifications

Arguably more interesting than the 1495 correctly classified sequences are the five incorrectly classified sequences (see Table I).

Of these, the misclassification of the insulated joint defect as a mechanical rail defect and the misclassification of the mechanical rail defect as an insulated joint defect are easily explained. The only difference between these sequences was the speed of the progression of the fault severity. This speed is drawn from normal distributions that are fault-dependent. Some realizations from these distributions will be very similar. In combination with the natural fluctuations of the current measurements, this will make some misclassifications inevitable.

The false negative misclassification of the ballast degradation sequences and the conductive object sequence seems to be related to the limits of the long-term time dependences that the network can handle.

For the misclassified conductive object sequence, the network inputs and outputs are shown in Fig. 5. The fault

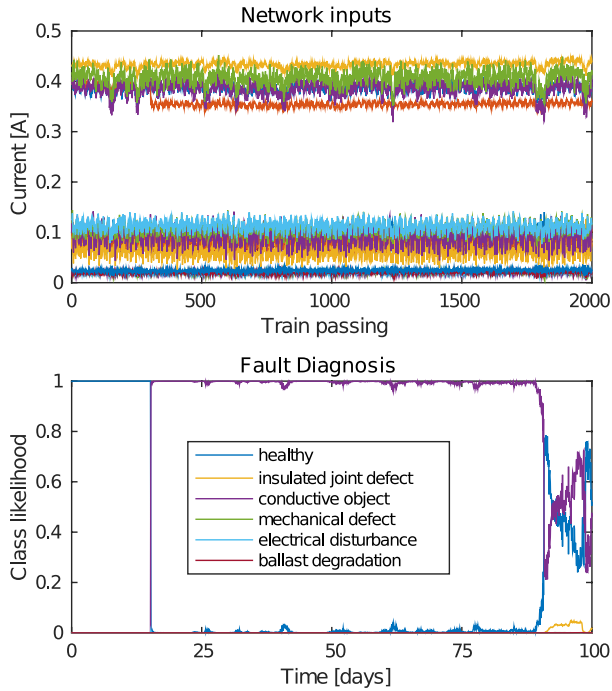


Fig. 5. Incorrectly classified sequence.

is initially classified correctly. Although this classification is kept for 1500 train passings, it seems that eventually new evidence makes the network forget the earlier events.

C. *t*-SNE

To gain some more insight into what the network has learned, the internal representations of the network at the end of the sequences will be investigated. After presenting each of the 1500 sequences to the network, the state of the memory units in the LSTM cells and the activations of the output units of the two recurrent layers in the network are stored. These activations are the network's internal representation of the sequence of events that has preceded the final time step and of the last input.

To compare these activation vectors, *t*-SNE [19] is used. This technique makes it possible to embed these 250-D vectors in a 2-D image in such a way that the vectors that are close together in the 250-D space are also close together in the 2-D plot. Therefore, sequences that are similar according to the network will occur close together in the plots. Note that the opposite does not have to be true; large distances do not necessarily imply that the sequences are very dissimilar.

1) *Role of the Layers*: The network has two hidden layers. The idea behind having multiple layers is that each subsequent layer uses the outputs of the previous layer to form higher level abstractions of the data. To investigate if this has happened, the activation vectors of the output units of both layers are plotted. Fig. 6 shows the activations of the output units in the first recurrent layer at the last time step for all 1500 sequences in the test set. Fig. 7 shows the same for the second layer.

From Fig. 6(a), it can be seen that the outputs of the first recurrent layer of the classification network are not

very sensitive to the temporal dependences in the data, as sequences from different classes are close together in the plot. From Fig. 6(b), it can be seen that the similarity of the outputs of the first layer seems to be based mostly on the fault severity at the final time step, as sequences with similar fault intensities are grouped close together.

The activation vectors of the output units in the second layer are labeled by the true fault category in Fig. 7(a). The grouping here seems based mostly on the true category and, therefore, on the underlying dependences that define these categories.

In Fig. 8(a), the state of the memory units in the second layer can be seen in the final time step of the sequences. It is interesting to note that the classes are less clearly separated here than they are in the output units of this layer. Presumably, the information about the fault severity coming from the first layer at the same time step is used to improve the classification. Alternatively, it might mean that the network remembers more information about the sequence than what is output at any given time to the softmax layer.

To gain more insights into how the network learns to classify faults, it can also be attempted to deduce how the network distinguishes between the conductive object and the electrical disturbance fault categories. Both faults abruptly lower the value of the current when a train is not present in the section. But where the current subsequently stays low for the conductive object fault, it is only intermittently low for the electrical disturbance. Furthermore, an electrical disturbance affects multiple track circuits along the same track, where a conductive object impacts only one. From Fig. 8(b), it can be seen that the network keeps a memory of a conductive object being present in the network. It does not, however, keep a memory of the fact that the electrical disturbances have been observed earlier in the sequence, as the sequences for which this is the case are not separated from those of the healthy sequences. In fact, in Fig. 8(a), it can be seen that also for the sequences that are at that time step undergoing an electrical disturbance, the state of the memory is similar to those in the healthy state.

2) *Spatial Dependences*: As discussed in Section III, the prior knowledge of the spatial and temporal fault dependences is not explicitly integrated into the network. Doing so on real data could introduce a bias if the prior knowledge turns out to be inaccurate. Since the neural network is trained and tested with synthetic data that are generated by a model that is based on the prior knowledge, it is interesting to see to what extent the network has learned to identify these dependences by itself.

Clearly, since the fault categories differ only based on their spatial and temporal dependences and the network manages to correctly classify them in 99.7% of the trials, it has learned to distinguish between these dependences. However, from Table III, it can be seen that the spatial dependences are not strictly necessary to distinguish between these five faults. Therefore, it is interesting to see if the network has learned these dependences or not.

The degradation of the ballast can affect either one track circuit or several along the same track. These spatial dependences are identified with D_1 and D_2 , respectively. For each sequence with a ballast degradation fault, one of these options

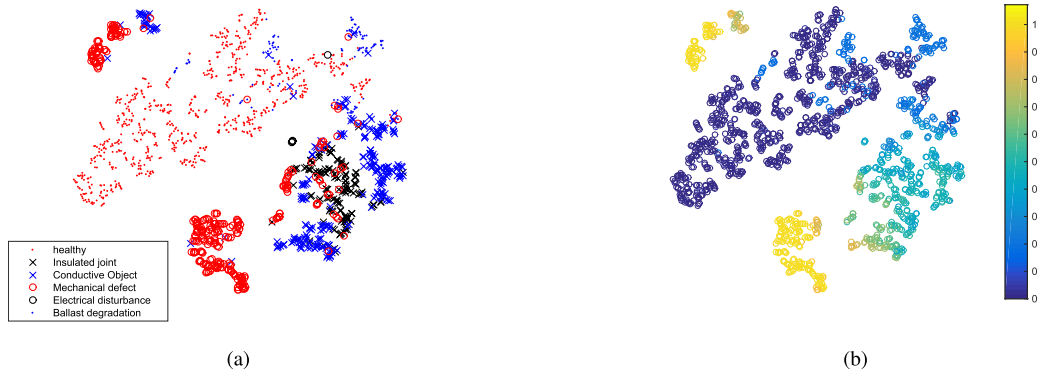


Fig. 6. t-SNE representation of the activity vectors of the output units in the first recurrent layer at the last time step of the sequences in the test data set $[h^1(2000)]$. (a) Labeled by true category. (b) Labeled by fault severity.

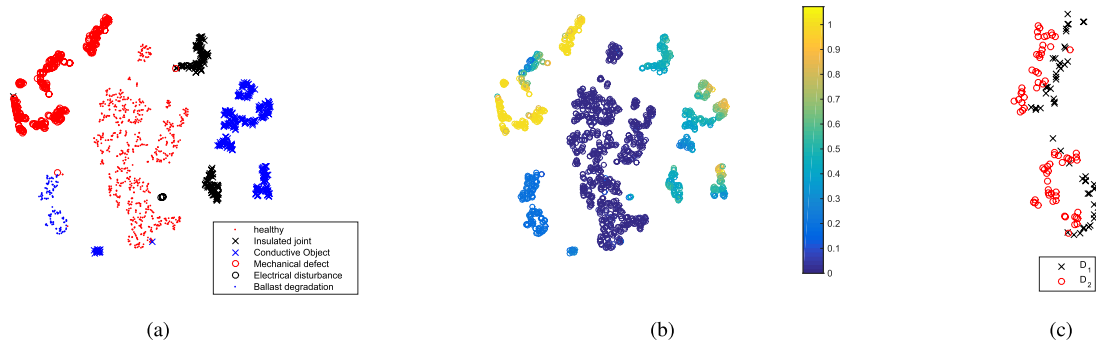


Fig. 7. t-SNE representation of the activity vectors of the output units in the second recurrent layer at the last time step of the sequences in the test data set $[h^2(2000)]$. (a) Labeled by true category. (b) Labeled by fault severity. (c) Ballast degradation sequences labeled per spatial dependence.

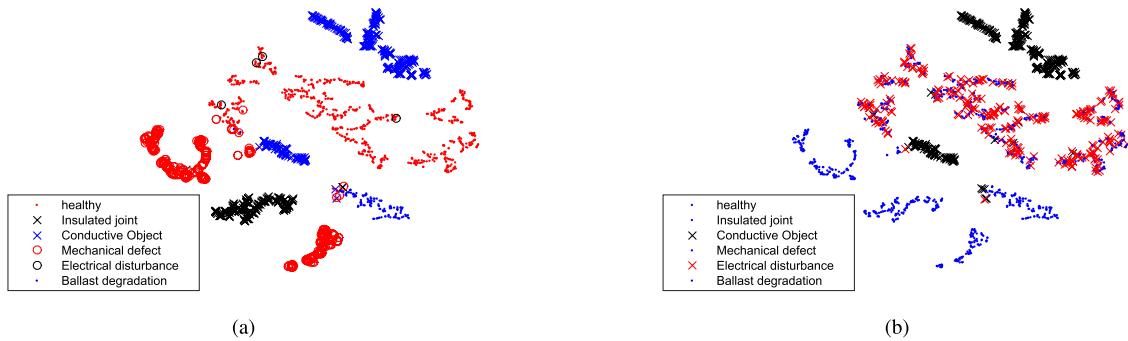


Fig. 8. t-SNE representation of the state of the memory units in the second recurrent layer at the last time step of the sequences in the test data set $[M^2(2000)]$. (a) Labeled by true category. (b) Separating conductive objects from electrical disturbances.

is picked with equal probability. In Fig. 7(c), the sequences suffering from the ballast degradation fault are shown. It appears from the plot that although these sequences are very similar, the network does distinguish between these spatial dependences.

V. CONVOLUTIONAL NETWORK COMPARISON

Besides LSTM-RNNs, convolutional neural networks (CNNs) [22] are a popular choice for dealing with temporal data [23]. In this section, we compare our LSTM network with a CNN.

The CNN that we consider is a feedforward network that takes all of the measurements of the past 2000 train passings

on the five track circuits at once as an input and gives the classification of the sequence at the most recent time step as an output. The CNN has two convolutional layers, followed by a fully connected layer with rectified linear unit nonlinearities and a softmax output layer. Both convolutional layers consist of two sublayers. The first performs a convolution step where a series of kernels is convolved with the inputs to the layer. The second sublayer performs a max-pooling step that takes the maximum activation of the kernels over a certain time window. The max-pooling operation introduces a limited invariance to the exact time at which a certain input pattern was detected. This simplifies the learning procedure and improves generalization. The kernel widths and the number

TABLE II

CONFUSION MATRIX OF THE CONVOLUTIONAL NETWORK FOR THE FAULT DIAGNOSIS TASK ON THE TEST DATA SET WITH 1500 SEQUENCES. THE ROWS INDICATE THE TRUE CLASS AND THE COLUMNS REPRESENT THE PREDICTED CLASS

true cat. / pred cat.	1	2	3	4	5	6
1 (healthy)	753	0	1	0	0	0
2 (Insulated joint defect)	0	132	0	0	0	0
3 (conductive object)	1	0	238	0	0	0
4 (mechanical rail defect)	5	1	0	245	0	0
5 (electrical disturbance)	4	0	0	0	0	0
6 (Ballast degradation)	0	0	0	0	0	121

of filters were chosen based on the prior knowledge of the faults and in such a way that the total number of parameters was approximately equal to that of the LSTM network.

Table II gives the classification results for the CNN. These can be compared with the results of our LSTM method, which are presented in Table I.

One thing that stands out is the inability of the CNN to diagnose electrical disturbances. This can be related to the max-pooling operator. This operator is relevant to most of the considered faults, as they degrade with time. So, if they were present at any previous time step, they are also present at the current time step. This is not the case for electrical disturbances, since these faults are intermittent and the classification of the sequence depends on the presence of the fault at the most recent time step. By introducing the time invariance through the max-pooling operation, the network is no longer applicable for diagnosing these faults. It has been found experimentally, however, that removing the max-pooling does not enable the convolutional network to correctly classify the electrical disturbances. The removal of the max-pooling step did result in reduced performance on the test data as the network started to overfit on the training data.

While the trained LSTM network is able to diagnose all types of faults with good accuracy, the learning performance is quite sensitive to the choice of the hyperparameters. The convolutional network gives slightly worse overall performance, but achieved this performance consistently for a wide range of hyperparameters, such as the kernel sizes, number of kernels, optimization algorithm, and learning rates. The training was also significantly faster.

In addition to the overall performance and ease of training, the suitability of the two methods differs per fault type. As discussed before, the LSTM network is more appropriate for intermittent faults and yields better overall performance. It does, however, sometimes forget faults that started a long time ago (see Section IV-B). Since the convolutional network does not use a memory, it does not suffer from this problem. Given the complementary strengths, it might be beneficial to combine both methods, as proposed in [24] and [25].

VI. CONCLUSION

In this paper, an LSTM-RNN has been proposed for fault diagnosis in railway track circuits. Synthetic data from a generative model are used to train and test the network.

TABLE III

FAULT TYPES AND THEIR SPATIAL AND TEMPORAL DEPENDENCES

Fault type	Spatial	Temporal	Fault rate
Insulated joint defect	D_1	$L \vee E$	intermediate
Conductive object	D_1	A	-
Mechanical rail defect	D_1	E	high
Electrical disturbance	D_2	I	-
Ballast degradation	$D_1 \vee D_2$	$L \vee E$	low

This enabled us to explore the opportunities of using this network in this setting. It has been shown that the network could learn the spatial and temporal dependences that characterize the considered faults directly from the electrical current measurements, without the manual integration of prior knowledge into the network. Of the 1500 scenarios presented to the network, 1495 were classified correctly. Furthermore, no false positive fault detections were made.

Although this paper has focused specifically on railway track circuits, LSTM-RNNs seem a promising option for other fault diagnosis problems as well, especially when the faults are characterized by long-term temporal dependences. We compared our LSTM network with a convolutional network. While the LSTM network outperforms the convolutional network for the track circuit case, the convolutional networks are easier to train. Given their complementary strengths, a combination of these networks might result in better performance on general fault diagnosis tasks than either of the individual networks can achieve.

APPENDIX A

FAULT TYPES CONSIDERED AND THEIR SPATIAL AND TEMPORAL DEPENDENCES

In this paper, we consider the following temporal dependences:

- L Linear;
- E Exponential;
- A Abrupt;
- I Intermittent.

Some faults that depend on time in a linear or exponential fashion can also be distinguished by the relative speed of the dependence. The spatial dependences considered are the following.

- D_c : The fault only affects the current in one track circuit.
- D_t : The fault affects the current in more track circuits on the same track.
- D_a : The fault affects the current of all track circuits in a certain area.

A. Fault Types Considered

In this paper, we consider a set of five different faults as described in the following. Table III gives a summary of the spatial and temporal dependences per fault type.

1) *Insulation Imperfections*: The sections of a railway track are electrically separated by insulated joints. When these joints wear out, the track circuit current of one section can leak into the adjacent section. The system is designed to be fail safe,

ensuring that the section that the current leaks into will not be identified as free because of this leakage. However, the current level in the section that the signal leaks out of will drop, potentially causing the section to be incorrectly identified as being occupied.

The effect of this fault will only be noticed in the section that the current leaks out of. As trains pass over the damaged joint, the defect will gradually get worse. The fault severity is, therefore, expected to increase either linearly or exponentially.

A conductive object placed over an insulated joint has a similar effect as the joint defect. In this case, however, the effect will occur abruptly and will not deteriorate over time.

2) *Rail Conductance Impairments*: The current travels through the rails from the transmitter to the receiver. When the impedance of this path increases, the current level in the receiver will decrease. One fault that can cause this problem is a mechanical defect in the rail. This fault would be specific to a single section and will increase exponentially over time as each passing train would cause greater damage to the deteriorating rail.

Another reason for the impedance of the rails to increase is the influence of disturbance currents. An example of this is when the track is saturated with traction currents. This problem occurs intermittently and affects several track circuits along the same track.

3) *Ballast Degradation*: Some current will always leak through the ballast between the rails in the section. The amount of current that leaks through the ballast depends on the impedance of the ballast. This impedance varies as a consequence of environmental conditions.

The ballast can also degrade over time, leading to a linear or exponential reduction in the magnitude of the signaling current when no train is present in the section. This effect would be noticeable in one or more sections along the same track. Compared with other faults, this fault would likely develop more slowly.

APPENDIX B GENERATIVE MODEL

To create a model that generates the amplitude of the electrical current $I(t)$ in the receiver of a track circuit as a train passes through the section, a data set of measurement sequences from $\mathcal{T} = 30\,000$ train passings has been studied. A mathematical model that was found to accurately describe these measurements was then fitted to the data. This model is based on four phases during a train-passing event.

- 1) *Phase 1*: Between t_0 and t_1 , the train has not yet arrived in the section. During this phase, the current $I(t)$ through the receiver should, therefore, be at the high level: $I(t) = I_h$.
- 2) *Phase 2*: At $t = t_1$, the first wheel set of the train enters the section. If the resistance of the wheel-set short circuit is low enough, this should result in a very quick drop of $I(t)$ to its low value I_l . However, in a large portion of the samples in the data set, the current drop is more gradual. By fitting a number of samples from three different track circuits to several equations for step responses, it was found that this phase could be

accurately and robustly described by an equation of the form $I(t) = \alpha_1 e^{-\tau_{\alpha 1}(t-t_1)} + \beta_1 e^{-\tau_{\beta 1}(t-t_1)}$.

- 3) *Phase 3*: Although ideally $I(t) = I_l$ should hold until the last wheel set of the train leaves the section, in the majority of the samples in the data set, the current starts to increase before this time. The curve between $t = t_2$ where the current is at the lowest level and $t = t_3$ where the last wheel set leaves the section can in almost all cases be accurately described by a function of the form: $I(t) = \alpha_2 e^{\tau_{\alpha 2}(t-t_2)} + \beta_2 e^{\tau_{\beta 2}(t-t_2)}$.
- 4) *Phase 4*: After the last wheel set leaves the section at $t = t_3$, the current $I(t)$ quickly increases to a value near I_h . On some of the samples, some overshoot is observed, and on some samples, a trend after the step is observed. Although a first-order step response was found to accurately describe many of the samples, a function of the form $I(t) = \alpha_3 e^{-\tau_{\alpha 3}(t-t_3)} + \beta_3 e^{\tau_{\beta 3}(t-t_3)}$ was found to represent these less common cases as well and is, therefore, chosen for the initial fitted model.

In Fig. 3(II), it can be seen that this model accurately describes the development of the current over time $I(t)$ during a train-passing event \mathcal{T} . This model was fitted to all of the measured data sequences. By analyzing the distributions of the values of the fitted model parameters, it was possible to create a simplified model with only a minimal sacrifice to the fitting accuracy. This model is given by

$$I(t) = I_l + \Delta I_{\max} \cdot \begin{cases} 1 & \text{for } t < t_1 \\ (1 - R)e^{-\tau_{\alpha 1}(t-t_1)} + Re^{-\tau_{\beta 1}(t-t_1)} & \text{for } t_1 \leq t < t_2 \\ (t - t_2) \frac{\Delta I_3}{(t_3 - t_2)} & \text{for } t_2 \leq t < t_3 \\ 1 - e^{-\tau_3(t-t_3)} & \text{for } t \geq t_3 \end{cases} \quad (9)$$

with the following values for the time constants

$$\begin{aligned} \tau_{\alpha 1} &= 9.25 \\ \tau_{\beta 1} &= 1.7 \\ \tau_3 &= 12.5 \end{aligned}$$

and

$$\Delta I_{\max} = I_h - I_l.$$

In this simplified model, the properties of the track circuit and the passing train are now represented by the four variables I_h , I_l , R , and ΔI_3 . By fitting the simplified model to the measured data and investigating the environmental conditions at the time of the measurements, the dependences of these four variables on several sources of normal variation were found. These sources include the precipitation, the time of day, and the train specific variations.

As these dependences only explain part of the observed variation in the measured data, several short- and long-term stochastic variations have been added to the model that affect both single track circuits and several track circuits in an area. In addition, the nominal parameters of the track circuits as well as the sensitivity of each track circuit to the sources

of variation are determined stochastically for each track circuit. This ensures that the synthetic data that the generative model produces contain comparable types of variation to the true measurement data. This makes it possible to not only determine the robustness of the condition monitoring method to these variations, but also its ability to pick up more subtle dependences in the data and use them, for example, whether the influences will affect all track circuits in a small area. By correctly identifying this influence, the effects on the measured signal could be filtered out, improving the condition monitoring performance.

B. Sampling Strategy

The faults that can affect the performance of the track circuits will in most cases change the values of the parameters I_h , I_l , R , and ΔI_3 very slowly over time. It is, therefore, important to sample the current $I(t)$ in a way that is informative of these values while taking as few samples per train passing \mathcal{T} as possible to ensure a high information density in the measurements.

Based on (9), the following sampling times are used.

- 1) t_1 : Just before the train arrives in the section: when the amplitude of the track circuit current is the highest I_h .
- 2) $t_1 + 0.35$ s: The value of the current $I(t)$ at 0.35 s after the first wheel set of the train enters the section is most instructive about the value of parameter R .
- 3) t_2 : When the current is at its lowest value I_l , about halfway through the train-passing event.
- 4) t_3 : Just before the last wheel set of the train leaves the section. This measurement gives ΔI_3 .

These four sampling times are shown in Fig. 3(II). As the simplified model of (9) fits the measured data well, these sampling times should also work for real measurement data.

To ensure a high enough information density for the artificial neural network to learn the long-term temporal fault dependences, these four current values are observed for all five considered track circuits and presented as one input time step \mathcal{T} to the network. This means that two trains (one on each track) should pass through the area before a new input is presented to the network.

REFERENCES

- [1] J. Chen, C. Roberts, and P. Weston, "Fault detection and diagnosis for railway track circuits using neuro-fuzzy systems," *Control Eng. Pract.*, vol. 16, no. 5, pp. 585–596, 2008.
- [2] K. Verbert, B. De Schutter, and R. Babuška, "Exploiting spatial and temporal dependencies to enhance fault diagnosis: Application to railway track circuits," in *Proc. Eur. Control Conf.*, Linz, Austria, Jul. 2015, pp. 3047–3052.
- [3] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [4] R. Wu, S. Yan, Y. Shan, Q. Dang, and G. Sun. (2015). "Deep image: Scaling up image recognition." [Online]. Available: <http://arxiv.org/abs/1501.02876>
- [5] A. Y. Hannun *et al.* (Dec. 2014). "Deep speech: Scaling up end-to-end speech recognition." [Online]. Available: <http://arxiv.org/abs/1412.5567>
- [6] L. Oukhellou, A. Debiolles, T. Dencœur, and P. Akin, "Fault diagnosis in railway track circuits using Dempster–Shafer classifier fusion," *Eng. Appl. Artif. Intell.*, vol. 23, no. 1, pp. 117–128, 2010.
- [7] Z. L. Cherfi, L. Oukhellou, E. Côme, T. Dencœur, and P. Akin, "Partially supervised independent factor analysis using soft labels elicited from multiple experts: Application to railway track circuit diagnosis," *Soft Comput.*, vol. 16, no. 5, pp. 741–754, 2012.
- [8] M. A. Sandizadeh and M. Dehghani, "Intelligent condition monitoring of railway signaling in train detection subsystems," *J. Intell. Fuzzy Syst., Appl. Eng. Technol.*, vol. 24, no. 4, pp. 859–869, 2013.
- [9] S. Sun and H. Zhao, "Fault diagnosis in railway track circuits using support vector machines," in *Proc. 12th Int. Conf. Mach. Learn. Appl.*, vol. 2. Miami, FL, USA, Dec. 2013, pp. 345–350.
- [10] Z. Lin-Hai, W. Jian-Ping, and R. Yi-Kui, "Fault diagnosis for track circuit using AOK-TFRs and AGA," *Control Eng. Pract.*, vol. 20, no. 12, pp. 1270–1280, 2012.
- [11] S. Ntalampiras, "Fault identification in distributed sensor networks based on universal probabilistic modeling," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 26, no. 9, pp. 1939–1949, Sep. 2015.
- [12] M. M. Gardner *et al.*, "Equipment fault detection using spatial signatures," *IEEE Trans. Compon., Packag., Manuf. Technol. C*, vol. 20, no. 4, pp. 295–304, Oct. 1997.
- [13] J. Chen, S. Kher, and A. Somani, "Distributed fault detection of wireless sensor networks," in *Proc. Workshop Dependability Issues Wireless Ad Hoc Netw. Sensor Netw.*, 2006, pp. 65–72.
- [14] G. E. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural Comput.*, vol. 18, no. 7, pp. 1527–1554, 2006.
- [15] J. Sun, R. Wyss, A. Steinecker, and P. Glocker, "Automated fault detection using deep belief networks for the quality inspection of electromotors," *tm-Technisches Messen (TM-TECH MESS)*, vol. 81, no. 5, pp. 255–263, 2014.
- [16] C. Shang, F. Yang, D. Huang, and W. Lyu, "Data-driven soft sensor development based on deep learning technique," *J. Process Control*, vol. 24, no. 3, pp. 223–233, 2014.
- [17] O. Obst. (2009). "Distributed fault detection in sensor networks using a recurrent neural network." [Online]. Available: <http://arxiv.org/abs/0906.4154>
- [18] H. C. Cho, J. Knowles, M. S. Fadali, and K. S. Lee, "Fault detection and isolation of induction motors using recurrent neural networks and dynamic Bayesian modeling," *IEEE Trans. Control Syst. Technol.*, vol. 18, no. 2, pp. 430–437, Mar. 2010.
- [19] L. van der Maaten and G. Hinton, "Visualizing data using t-SNE," *J. Mach. Learn. Res.*, vol. 9, nos. 2579–2605, p. 85, 2008.
- [20] A. Graves and N. Jaitly, "Towards end-to-end speech recognition with recurrent neural networks," in *Proc. 31st Int. Conf. Mach. Learn. (ICML)*, 2014, pp. 1764–1772.
- [21] P. J. Werbos, "Backpropagation through time: What it does and how to do it," *Proc. IEEE*, vol. 78, no. 10, pp. 1550–1560, Oct. 1990.
- [22] Y. LeCun and Y. Bengio, "Convolutional networks for images, speech, and time series," *The Handbook of Brain Theory and Neural Networks*. Cambridge, MA, USA: MIT Press, 1998.
- [23] M. Långkvist, L. Karlsson, and A. Loutfi, "A review of unsupervised feature learning and deep learning for time-series modeling," *Pattern Recognit. Lett.*, vol. 42, pp. 11–24, Jun. 2014.
- [24] T. N. Sainath, O. Vinyals, A. Senior, and H. Sak, "Convolutional, long short-term memory, fully connected deep neural networks," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, Apr. 2015, pp. 4580–4584.
- [25] L. Deng and J. C. Platt, "Ensemble deep learning for speech recognition," in *Proc. INTERSPEECH*, 2014, pp. 1915–1919.



Tim de Bruin received the B.Sc. degree in mechanical engineering and the M.Sc. degree in systems and control from the Delft University of Technology, Delft, The Netherlands, in 2012 and 2015, respectively, where he is currently pursuing the Ph.D. degree with the Delft Center for Systems and Control.

His current research interests include neural networks, reinforcement learning, and robotics.



Kim Verbert received the B.Eng. (*cum laude*) degree in human kinetic technology from the Hague University of Applied Sciences, The Hague, The Netherlands, in 2009, and the M.Sc. (*cum laude*) degree in control engineering from the Delft University of Technology, Delft, The Netherlands, in 2012, where she is currently pursuing the Ph.D. degree with the Delft Center for Systems and Control.

Her current research interests include fault diagnosis, maintenance optimization, friction compensation, and human motion control.



Robert Babuška received the M.Sc. (Hons.) degree in control engineering from Czech Technical University in Prague, in 1990, and the Ph.D. (*cum laude*) degree from the Delft University of Technology (TU Delft), the Netherlands, in 1997. He has had faculty appointments with the Czech Technical University in Prague and with the Electrical Engineering Faculty, TU Delft. Currently, he is a Professor of Intelligent Control and Robotics at the Delft Center for Systems and Control.

He is also the director of the TU Delft Robotics Institute. His current research interests include reinforcement learning, neural and fuzzy systems, nonlinear identification, state-estimation, model-based and adaptive control and dynamic multi-agent systems. He has been involved in the applications of these techniques in the fields of robotics, mechatronics, and aerospace.