

Detecting Message Modification Attacks on the CAN Bus with Temporal Convolutional Networks

Irina Chiscop¹^a, András Gazdag²^b, Joost Bosman¹^c and Gergely Biczók²^d

¹*Cyber Security & Robustness Department, TNO, The Hague, The Netherlands*

²*CrySyS Lab., Dept. of Networked Systems and Services, Budapest University of Technology and Economics, Budapest, Hungary*

Keywords: Vehicle Security, Intrusion Detection, Controller Area Network, Machine Learning, Temporal Convolutional Networks.


Abstract: Multiple attacks have shown that in-vehicle networks have vulnerabilities which can be exploited. Securing the Controller Area Network (CAN) for modern vehicles has become a necessary task for car manufacturers. Some attacks inject potentially large amount of fake messages into the CAN network; however, such attacks are relatively easy to detect. In more sophisticated attacks, the original messages are modified, making the detection a more complex problem. In this paper, we present a novel machine learning based intrusion detection method for CAN networks. We focus on detecting message modification attacks, which do not change the timing patterns of communications. Our proposed temporal convolutional network-based solution can learn the normal behavior of CAN signals and differentiate them from malicious ones. The method is evaluated on multiple CAN-bus message IDs from two public datasets including different types of attacks. Performance results show that our lightweight approach compares favorably to the state-of-the-art unsupervised learning approach, achieving similar or better accuracy for a wide range of scenarios with a significantly lower false positive rate.


1 INTRODUCTION


A modern automobile contains more than 100 electronic control units (ECU) to control the vehicular subsystems and help the driver with various sophisticated services. These ECUs are spread over the entire vehicle and connected mostly via Controller Area Networks (CANs). CAN was initially designed to be an isolated system decades ago, however, in the age of round-the-clock connected networked objects, this property no longer holds. Realizing the aforementioned connectivity, Bluetooth, Wi-Fi, or cellular connections are all potential intrusion points for an attacker. If a malicious actor compromises a component^{1,2} that implements one of these connections, it becomes possible to manipulate the CAN net-


work. Dashboard information is displayed and ctuators in Advanced Driver-Assistance Systems (ADAS) are controlled based on sensor readings transmitted over the CAN bus; interfering with these messages may result in significant financial loss and danger to human life.

Two approaches have so far been seen among the attacks. Existing attacks are of two distinct types: (i) the attacker either injects additional CAN messages into the network or (ii) she modifies otherwise valid messages sent by legitimate ECUs. We note that the latter attack is very difficult to implement by altering electrical signals of the CAN on-the-fly, yet it can be realized by compromising an ECU and sending out messages with modified content, or by compromising a gateway between two CAN networks and modifying messages passing through it (Gazdag et al., 2020). Message injection attacks can be detected easily because the original messages are also present next to malicious ones, which changes the temporal patterns of traffic. On the other hand, detecting modification attacks poses a tougher challenge: it requires an in-depth analysis of the actual payload as the rest of the

^a <https://orcid.org/0000-0002-1249-8518>

^b <https://orcid.org/0000-0002-4481-3308>

^c <https://orcid.org/0000-0001-6325-1462>

^d <https://orcid.org/0000-0002-3891-3855>

¹www.wired.com/2015/07/hackers-remotely-kill-jeep-high-way

²www.wired.com/story/tesla-model-x-hack-bluetooth

traffic properties remain intact in this scenario.

Related Work. In recent years, a considerable amount of literature has been published on CAN bus intrusion detection. These works can be split into three categories: frequency-, statistics-, or machine learning based methods. Most of these approaches are particularly useful for detecting cyber-attacks in which additional messages are being injected into the CAN bus. The simplest of the three, frequency-based models focus on testing inter-arrival times of CAN messages against a predefined normal baseline (Taylor et al., 2015; Song et al., 2016; Moore et al., 2017; Gazdag et al., 2019). As the name suggests, statistics-based detection approaches exploit the statistical properties of CAN bus traffic such as entropy (Muter and Asaj, 2011), Z-score (Tomlinson et al., 2018) or Mahalanobis distance (Mo et al., 2019). Machine learning based methods imply the usage of artificial neural networks, clustering and supervised models for classification and regression. In the specific field of CAN bus intrusion detection, popular machine learning approaches include autoencoders (Lokman et al., 2019; Lin et al., 2021; Novikova et al., 2020), recurrent neural networks (RNN) such as Long Short-Term Memory (LSTM) networks (Taylor et al., 2016; Negi et al., 2019; Khan et al., 2020; Hanselmann et al., 2020; Hossain et al., 2020), Gated Recurrent Unit (GRU)-based networks (Kukkala et al., 2020), replicator neural networks (Weber et al., 2018), and deep convolutional networks (Song et al., 2020). The scrutinized literature shows that recurrent architectures are often the preferred choice for modeling the time series of CAN bus signals, whilst convolutional networks are used when data is transformed to a two-dimensional grid dataframe to resemble an image format (Song et al., 2020). In particular, only one approach was found to combine these two techniques in the form of a convolutional LSTM (Tariq et al., 2020) which is trained on labeled data in a supervised fashion.

Temporal Convolutional Networks. To the best of our knowledge, no existing solution employs (causal) convolutions to model the time series representation of CAN signals; we argue that such an approach makes perfect sense given the successful application of convolutional networks to sequence modeling tasks. Specifically, a Temporal Convolutional Network (TCN) is a type of convolutional network whose architecture consists of causal (and dilated) convolutions (Bai et al., 2018). It has been shown that this new type of network outperforms recurrent architectures, such as LSTM and GRU, on a multitude of sequence modeling tasks including the adding problem and image classification on sequential MNIST

and P-MNIST (Bai et al., 2018). In fact, TCNs have also been successfully applied to anomaly detection in general time series data (He and Zhao, 2019).

Our Contribution. In this paper, we propose a TCN-based approach for detecting modified CAN bus messages; our focus is solely on message modification attacks with no message injection. We construct and train the TCN in an unsupervised fashion, since, in practice, labelling CAN bus messages is a very difficult task. In the training process, the TCN will learn to accurately reconstruct the signals of individual CAN bus messages through its causal convolution layers, which allows for information retention from past data samples. Finally, the classification of new data samples will resume to setting an appropriate threshold on their reconstruction loss value. The core idea here is that signals whose data have been altered will be poorly reconstructed by the model, and thus be easy to recognize. Note, that it is not a prerequisite for us to know CAN bus signal semantics which varies for vehicle make and model, and is usually kept confidential (Lestyan et al., 2019; Remeli et al., 2019). The contribution of this paper is three-fold:

1. We first introduce a new CAN bus dataset containing both benign data and synthetic attacks.
2. We then propose a TCN architecture to learn and reconstruct the normal behaviour of CAN bus signals, and use this information to pinpoint anomalies that do not conform to the reconstruction given by model.
3. We compare the detection performance of our approach to a state-of-the-art GRU-autoencoder (Kukkala et al., 2020) (shown to outperform other existing solutions) through numerical experiments on both our own dataset and the *de facto* standard SynCAN dataset (Hanselmann et al., 2020). Results show that our simple TCN-based approach compares favorably to the state-of-the-art, i.e., it achieves similar or better accuracy with a significantly lower false positive rate.

Paper Structure. The rest of the paper is structured as follows. Section 2 presents our proposed TCN architecture in detail. Section 3 describes the design of our experiments including choosing the baseline, introducing our two datasets and the training process, and defining evaluation metrics. Section 4 presents the results of the comparative performance evaluation. Finally, Section 5 concludes the paper.

2 INTRUSION DETECTION MODEL

In this section we present the motivation behind choosing temporal convolutional networks as an intrusion detection mechanism for the CAN bus. We first provide some background on convolutional networks and then describe our proposed TCN architecture in detail.

2.1 Convolutional Networks

Convolutional neural networks are a particular kind of deep neural networks that enables the extraction of relevant spatial and temporal features from the input (e.g., an image) by learning a set of filters. These filters represent multi-dimensional arrays sliding over the input image, and are initialized randomly. During the forward pass, the dot product between the entries of each filter and the image sub-block is computed, resulting in a feature map. When another convolutional layer is added, the features learned in the first layer are combined to create new ones. To account for as many (non-linear) combinations of features as possible, it is customary to increase the filter size in the subsequent layers. The deeper the network becomes, the better it gets at extracting refined patterns from the data. A more detailed description of different convolutional architectures can be found in (Aloysius and Geetha, 2017).

Temporal convolutional networks (TCN) are a category of convolutional networks particularly suitable for modeling long-term dependencies in sequential data (van den Oord et al., 2016; Bai et al., 2018). Consider for instance the following task: based on input sequence x_0, x_1, \dots, x_T , predict corresponding output y_0, y_1, \dots, y_T at each time step. There are two constraints associated with this task. First, the predicted output y_t should only be influenced by previously observed inputs x_0, x_1, \dots, x_t , and, second, the size of the network output must be identical to that of the input sequence. TCNs tackle the first constraint by sliding a filter only over the past input values. In other words, the convolution filter has positive weights only for past inputs. TCNs also employ dilated causal convolutions which, unlike regular causal convolutions, enable an exponential growth of the receptive field by skipping over the inputs while convolving. Moreover, a larger receptive field allows the neural network to infer the relationships between different observations in the input data. The second constraint is addressed by padding the input data with zeros at the borders, to control the dimension of the output. These two architectural elements can be observed in Figure 1, depict-

ing a dilated causal convolutional network with two hidden layers. Here, the zero-padding is represented by the white squares on the left side. The filter size of $k = 3$ is indicated by the blue lines. The dilation factor d , applied at each layer, indicates how many input values are being skipped by the filter. Increasing the dilation factor by 2 at each subsequent layer results in a receptive field of size 15: the value of a neuron in the output layer is influenced by fifteen neurons from the input layer.

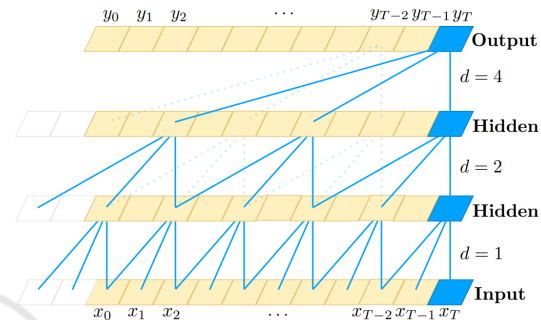


Figure 1: A dilated causal convolutional neural network with two hidden layers, dilation factors $d = 1, 2, 4$ and filter size $k = 3$ (Bai et al., 2018).

TCNs possess numerous advantages when compared to recurrent architectures (Bai et al., 2018). Convolutions within TCNs can be computed in parallel, thus allowing the entire data sequence to be processed. That is not possible with RNNs, where the computation of the output at a specific timestep requires the complete computation of all its predecessors. Moreover, TCNs require less memory during training than RNNs, where partial values of cell-gates need to be stored, and exhibit stable gradients, as backpropagation does not happen through multiple different time samples. In theory the receptive field of RNNs is infinite; in TCNs the field is finite, and its size depends on the number of layers (dilations) and filters used. Apparently, there exists a trade-off between how lightweight the network is, and its ability to capture long-term dependencies in the data. Both aspects are equally important to obtain a scalable and reliable CAN bus intrusion detector. In the remainder of this paper we show that a TCN model is a suitable candidate for this purpose.

2.2 TCN Architecture

The TCN to be used for CAN bus intrusion detection follows the general framework from (Bai et al., 2018) and is shown in Figure 2. The network consists of an input layer, three residual blocks, and an output layer. As shown in the figure, the input for the TCN must be

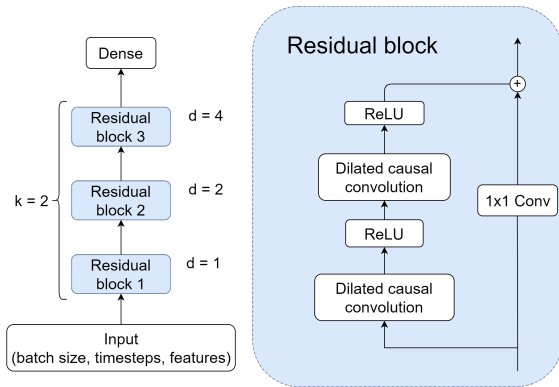


Figure 2: Our TCN architecture with three residual blocks with convolutional dilations and filter size of $k = 2$.

three-dimensional. Each residual block contains two dilated causal convolution layers each having 64 filters and the same dilation factor d . The Rectified Linear Unit (ReLU) is used as an activation function on these layers. The filter size is kept at the same value of $k = 2$ across all residual blocks. A skip connection is also enabled, which adds the output from the previous layer to the next layer. This is marked by the element-wise addition \oplus . Due to zero-padding, this operation may receive inputs that differ in shape. To circumvent this, a 1×1 convolution is added.

The network is kept simple deliberately: no weight normalization or dropout layers have been used. Our main objective here is to investigate whether this lightweight TCN can successfully learn to reconstruct CAN bus signals, and achieve results comparable to or better than other, more complex state-of-the-art classifiers.

2.3 Intrusion Score and Output

We distinguish between benign and malicious messages by applying a threshold to the reconstruction loss. We therefore monitor the squared error between the signal value at a given times and its latest reconstructed value. This defines an intrusion score for each signal in a message. To compute an intrusion score per message, we calculate a set of thresholds given by the 99.9th percentile of the validation loss for each signal in the data. A message is then labeled as malicious if one of the signal’s intrusion scores exceeds the threshold set for that signal. We opted for this approach to label messages based on individual signal thresholds: in practice, depending on the complexity and correlation of the signals, some may be better reconstructed during training than others.

3 EXPERIMENT DESIGN

In this section we describe the design of our numerical experiments, including the baseline model, datasets, training process and our choice of evaluation metrics.

3.1 Selecting the Most Suitable Baseline

For evaluation purposes, we identified the best-performing CAN bus anomaly detection algorithms by scrutinizing recent literature. We used the following selection criteria:

- Unsupervised learning: the algorithm requires no labeled data for training.
- Generalization: the algorithm is easy to generalize, and thus does not depend on data pre-processing such as identifying and pre-selecting specific CAN signals.
- Fully-reproducible: the algorithm needs to be accompanied by sufficient information in order to have a fully reproducible implementation.

To the best of our knowledge, the most recent and suitable candidate is the INDRA framework (Kukkala et al., 2020). It proposes a recurrent autoencoder network that is able to detect CAN messages in which signals have been tampered with. For each message ID one such recurrent autoencoder is trained such that it learns to reconstruct the signals within that particular message ID. This approach is shown to outperform other recent unsupervised methods such as Predictor LSTM (Taylor et al., 2015), Replicator Neural Network (Weber et al., 2018), and CANet (Hanselmann et al., 2020), on most attack classes of the SynCAN dataset, in terms of accuracy and false positive rate. Moreover, Predictor LSTM is designed to predict the raw message data in string form, and thus does not directly fall within the scope of time-series-based intrusion detection. Note that the CANet model is also more complex since its architecture combines the LSTM models of individual messages to account for capturing the correlations between different IDs. Finally, the convolutional LSTM proposed in (Tariq et al., 2020) is a promising method for predicting multi variate time series data. However, it was designed for supervised learning which requires labeled data for training and for this reason, it falls outside the scope of this paper.

In view of these arguments, INDRA is the most sensible baseline for comparative performance evaluation.

3.2 Datasets

When evaluating machine learning classifiers, it is considered best practice to employ multiple datasets in order to assess the impact of the number of data samples and different features on the model's performance. Moreover, publicly available CAN bus datasets for intrusion detection are labelled differently, either per message ID or per signal. To account for both, we consider two datasets: the SynCAN dataset with message labels and the CrySyS dataset with individual signal labels.

3.2.1 SynCAN Dataset

The SynCAN (Synthetic CAN Bus Data) dataset was introduced in (Hanselmann et al., 2020), and is publicly available³. The dataset contains 10 different CAN message IDs, whilst the number of signals in each ID varies between 1 and 4. Overall, the dataset covers 20 correlated signals. The training data spans approximately 16.5 hours of traffic, while the testing data about 7.5. Moreover, testing data includes a 0/1 label per individual message, to indicate whether it is malicious or not. However, there is no indication as to which signal has been attacked within a malicious message. Since this dataset is only meant for unsupervised learning purposes, the training data does not include explicit labels. Finally, the test data is split across six different files, each corresponding to a different simulated attack:

- *Plateau Attack*: the value of a single signal is overwritten by a constant value over a certain period of time.
- *Continuous Change Attack*: the value of a signal is overwritten at a slow pace, such that it increasingly deviates from its true value.
- *Playback Attack*: the values of a signal within a time interval is overwritten with the values of the same signal from a randomly selected past interval.
- *Suppression Attack*: signal values contained in a certain message ID simply do not appear in the CAN traffic for a period of time.
- *Flooding Attack*: messages with a certain ID are sent with a higher frequency to the CAN bus.

Detection of message injection attacks (suppression attack and flooding attack) is not a goal of this paper. Nonetheless, in Section 4, we evaluated our TCN architectures performance on those as well for a better comparison with the INDRA model.

³www.github.com/etas/SynCAN

3.2.2 CrySyS Dataset

The CrySyS dataset was created by the CrySyS Lab in the context of the SECREDAS project⁴, and it is also publicly available⁵. It is significantly smaller compared to the SynCAN dataset, however, the driving environment and the behavior of the vehicle are better known. It contains 7 smaller (<1 minute) captures of specific driving and traffic scenarios, and a longer trace (~25 minutes). There are 20 different message IDs in the traces, and the number of signals varies between 1 and 6.

Additionally, to complement this dataset, we have developed a signal extractor and an attack generator script. The signal extractor is based on the work presented in (Stone et al., 2018). It calculates the statistical properties of bits in CAN data fields, and identifies and separates the signals based on the changes in these values.

The attack generator⁶ is able to modify the CAN messages. It changes some or all the values in the data field of existing messages without modifying the timing of a message. To achieve a meaningful targeted attack, the generator can be combined with the information gathered from the signal extractor to modify specific signals in the trace. The attack generator also supports multiple attack types:

- *Change to Constant*: the original value is replaced by the given constant value.
- *Change to Random*: the original value is replaced by a new random value.
- *Modify with Delta*: the given value is added to the original data value.
- *Modify with Increment*: a per message increment is added to the original value.
- *Modify with Decrement Value*: a per message decrement is subtracted from the original value.
- *Change to Increment*: the original data value is replaced by a per message incremented value.
- *Change to Decrement*: the original data value is replaced by a per message decremented value.

We modified the original CrySyS traces with the attack generator script to simulate attacks. After we identified the different signals in the traces we replaced a chosen signal with a constant value for the second half of the trace. Note that this simple change-to-constant/plateau attack was enough to demonstrate the capabilities of our approach over INDRA (see

⁴www.secredas-project.eu

⁵www.crysys.hu/research/vehicle-security

⁶www.github.com/CrySyS/can-log-infector

Table 1: Overview of datasets used in the numerical experiments.

Dataset	Message ID	No. of signals	Train samples	Test samples
SynCAN	2	3	4139826	909869
	3	2	2070144	1884235
	10	4	1380087	610294
CrySyS	280	4	157472	3895
	290	5	15748	389

Section 4). Also note that we focused on IDs with 1 to 4 signals per message, similar to SynCAN, to be able to compare the results across the two datasets.

3.3 Training the Models

Training both the TCN and INDRA models required the normalization of signal data (values between 0 and 1), and then re-shaping the input data to three-dimensional. This was done by sliding a fixed-size window over the time series, one timestamp at a time. As in (Kukkala et al., 2020), we applied a rolling window of 20 timestamps or, equivalently, of 20 messages, to the training datasets shown in Table 1.

The rest of the training parameters were set to the same values as in (Kukkala et al., 2020) to ensure an accurate reproduction of the INDRA model. Concerning the optimizer and loss function, both models used the *Adam* optimizer with learning rate 0.0001 and mean square error. The models were trained for 100 epochs with a batch size of 128 on 85% of the training data, whilst the other 15% was kept for validation. An early-stop mechanism terminated the training if the validation loss did not improve in the last 10 epochs. Note that during initial experiments, a higher number of epochs was considered, but the training stopped before the 100th epoch in all cases. All models were implemented using the *keras* and *keras-tcn*⁷ libraries in Python 3.7, and trained on a GeForce GTX 960 GPU. The two models have only been trained offline, not on live CAN bus data.

3.4 Evaluation Metrics

To evaluate the performance of the TCN model, we use the intrusion score defined in Section 2.3. The INDRA model uses the same squared error as a signal intrusion score, but applies a generic threshold set to the 99.9th percentile of the validation loss (computed across all signals). The message intrusion score is then given by the maximum signal intrusion score contained in that message, and is then compared to the threshold. We use three standard performance metrics

⁷www.github.com/philipperemy/keras-tcn

for the evaluation of the models: accuracy, false positive rate and precision. Accuracy measures the ratio of the predicted labels exactly matching the ground truth, and is defined as follows:

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}, \quad (1)$$

where TN , TP , FN , FP denote the number of true, and false, positives and negatives, respectively. Accuracy gives an indication of the general classification capabilities of a certain model.

The false positive rate (FPR) measures the amount of samples wrongly classified as malicious, whilst in fact being benign. The false positive rate is extremely relevant from the practical point of view: in the CAN bus context, the messages marked as malicious may need to be further analyzed before deciding on mitigation actions. To keep operation efficient, the false positive rate needs to be minimized as much as possible. Precision, on the other hand, measures the capabilities of the model to actually detect the relevant attacks (positive samples). This is another important quantity to monitor since imbalanced datasets, with far more negatives than positives, may render accuracy a deceiving metric. In fact, CAN bus datasets are usually imbalanced, since most (simulated) attacks have a very short duration. The FPR and precision are defined as follows:

$$FPR = \frac{FP}{TN + FP}, \quad (2)$$

$$Precision = \frac{TP}{TP + FP}. \quad (3)$$

4 RESULTS

SynCAN. We first assessed the performance of the two models on the SynCAN dataset. The accuracy and false positive rate, calculated for the normal test set and for each attack class, are shown in Table 2. A first observation is that TCN achieves a higher accuracy than INDRA in most cases, with the exception of playback and flooding attacks on ID 10. Moreover, the false positive rates are quite low for both models, which can be explained by looking at the precision values in Table 3. Overall, there are large variations in the precision values across different message IDs which may be related to how the attacks were performed (target signals chosen, attack duration, etc.) and the different signal correlations. Also, the relatively low precision values in Table 3 show that the models manage to capture only a limited set of temporal characteristics of the SynCAN data. This is a

Table 2: Results for the SynCAN dataset.

Model	Data	Normal		Continuous		Playback		Flooding		Suppress		Plateau	
		Acc.	FPR	Acc.	FPR	Acc.	FPR	Acc.	FPR	Acc.	FPR	Acc.	FPR
TCN	ID 2	0.9977	0.0022	0.8660	0.0018	0.8674	0.0013	0.7678	0.0026	0.8402	0.0001	0.8336	0.0066
INDRA		0.9811	0.0188	0.8584	0.0121	0.8660	0.0046	0.7600	0.0157	0.8347	0.0101	0.8133	0.0495
TCN	ID 3	0.9992	0.0007	0.8664	0.0009	0.8680	0.0002	0.6422	0.0011	0.8390	0.0004	0.8394	0.0012
INDRA		0.9965	0.0034	0.8653	0.0033	0.8672	0.0012	0.6420	0.0033	0.8377	0.0025	0.8386	0.0036
TCN	ID 10	0.9977	0.0022	0.8637	0.0072	0.8577	0.0160	0.7399	0.0001	0.8446	0.0011	0.8282	0.0136
INDRA		0.9858	0.0141	0.8546	0.0176	0.8638	0.0070	0.7923	0.0047	0.8370	0.0105	0.8100	0.0447

Table 3: Precision of the models on SynCAN dataset.

Model	Data	Continuous	Playback	Flooding	Suppress	Plateau
TCN	ID 2	0.4457	0.2458	0.0205	0.9027	0.3022
INDRA		0.1992	0.3696	0.1577	0.2812	0.3033
TCN	ID 3	0.5231	0.0000	0.1028	0.5854	0.7809
INDRA		0.0143	0.0000	0.3766	0.3261	0.6192
TCN	ID 10	0.3706	0.1949	0.000	0.0212	0.2036
INDRA		0.1779	0.1668	0.9413	0.0386	0.2224

direct consequence of the stopping mechanism implemented during training, and in the case of TCN, of the choices made to keep a lightweight architecture. For playback attacks, precision is very low for both models, which leads to the similarly low false positive rates achieved in this class. This is not surprising since during a playback attack, a portion of past data is written over its current values, making the signal look normal, and thus the attack difficult to detect. TCN clearly achieves a better performance than INDRA in detecting continuous attacks. Moreover, for message IDs 2 and 3, TCN detects suppression attacks with a much larger precision compared to INDRA. This result appears to be influenced by the number of signals in the message, since precision significantly decreases as the number of signals increases. As for plateau attacks, the two methods achieve similar results. INDRA is more precise than the TCN model is on detecting flooding attacks. This is an expected result, mainly due to the TCN accurately reconstructing data from a flooding attack since the data values are not altered during such an attack. To sum up, the TCN model is capable of detecting all message modification attacks (continuous change, playback and plateau) effectively. Although detecting attacks which modify the arrival rates of CAN bus messages was not part of the original goal, TCN also proved successful at detecting suppression attacks.

CrySyS. The message IDs in the SynCAN dataset contains signals that are physically interdependent, but are very weakly correlated; this also increases the difficulty of the detection task. In order to assess how the two models perform in a different setting, we consider two message IDs of the CrySyS dataset which contains more signals with a strong correlation. Here, similarly to SynCAN, only one signal was attacked. The results are shown in Table 4. We notice that both

Table 4: Results for the CrySyS dataset.

Model	Data	Acc.	FPR	Precision
TCN	ID 280	0.8833	0.0426	0.7766
INDRA		0.7989	0.0000	0.0000
TCN	ID 290	0.9159	0.0687	0.7701
INDRA		0.8617	0.0378	0.7755

models still achieve high accuracy and a low false positive rate, with TCN showing a high precision for both attacks, as opposed to INDRA, failing to detect the attack in message 280.

Summary. The simple TCN architecture achieves a slightly better accuracy compared to the INDRA model on both datasets. A remarkable achievement of TCN is the significant reduction of false positives (by a factor of 10) in nearly all cases: this translates to a more reliable detector in practice. Another advantage of the TCN is that it is quick to train, and achieves in general lower training and validation loss (see Figure 3 for an example).

5 CONCLUSIONS

In this paper we examined the applicability of temporal convolutional networks to CAN bus intrusion detection, with a focus on message modification attacks. To this end, we proposed a lightweight TCN, and showed that its classification performance compared favorably to the state-of-the-art baseline INDRA across different datasets and attack classes. Specifically, we demonstrated that our computation-efficient and compact TCN model achieves similar or better accuracy, while reducing false positives with an order of magnitude. This shows that TCNs have a great potential both in modeling CAN bus signal and being deployed in practical settings.

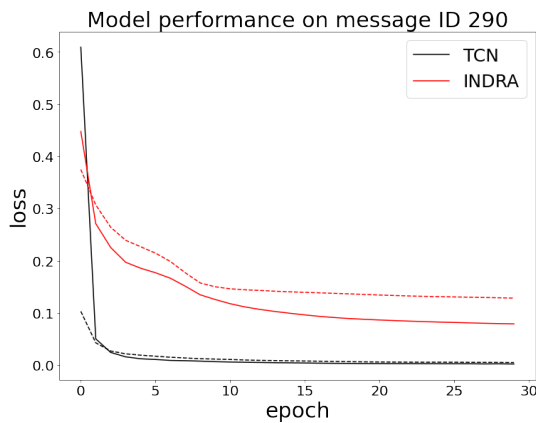


Figure 3: Training loss (continuous lines) and validation loss (dashed lines) of the two models on message ID 290 of the CrySyS dataset.

Future Work. First of all, the elimination of the early termination mechanism would potentially yield better performance; early termination was necessary in our experiments due to hardware-related constraints. Second, the TCN architecture was kept very simple on purpose to ensure a computationally lightweight model. However, the learning abilities of the network could be improved by increasing the filter size and the dilation factor between causal convolutions, and by stacking additional residual blocks together. Third, it is worthwhile to investigate how message-based and signal-based intrusion thresholds, and the underlying intra-message signal correlation influence the performance of both models for different attack classes. Finally, correlations between signals across different message IDs could be considered leading to a more accurate representation of normal CAN bus behaviour. To this end, an architecture combining multiple TCN blocks (modeling individual message IDs as a la CANet (Hanselmann et al., 2020)) could be used.

ACKNOWLEDGEMENTS

This work has been partially funded by the European Commission via the H2020-ECSEL-2017 project SECREDAS (Grant Agreement no. 783119). The research presented in this paper and carried out at the Budapest University of Technology and Economics have been supported by the NRD Office, Ministry of Innovation and Technology, Hungary, within the framework of the Artificial Intelligence National Laboratory Programme, and the NRD Fund based on the charter of bolster issued by the NRD Office.

REFERENCES

- Aloysius, N. and Geetha, M. (2017). A review on deep convolutional neural networks. In *2017 International Conference on Communication and Signal Processing (ICCSP)*, pages 0588–0592.
- Bai, S., Kolter, J. Z., and Koltun, V. (2018). An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *ArXiv*, abs/1803.01271.
- Gazdag, A., Ferenczi, C., and Buttyán, L. (2020). Development of a man-in-the-middle attack device for the can bus. In *1st Conference on Information Technology and Data Science*.
- Gazdag, A., Neubrandt, D., Buttyán, L., and Szalay, Z. (2019). Detection of injection attacks in compressed can traffic logs. In Hamid, B., Gallina, B., Shabtai, A., Elovici, Y., and Garcia-Alfaro, J., editors, *Security and Safety Interplay of Intelligent Software Systems*, pages 111–124, Cham. Springer International Publishing.
- Hanselmann, M., Strauss, T., Dormann, K., and Ulmer, H. (2020). CANet: An unsupervised intrusion detection system for high dimensional can bus data. *IEEE Access*, 8:58194–58205.
- He, Y. and Zhao, J. (2019). Temporal convolutional networks for anomaly detection in time series. *Journal of Physics: Conference Series*, 1213:042050.
- Hossain, M. A., Inoue, H., Ochiai, H., Fall, D., and Kadobayashi, Y. (2020). LSTM-based intrusion detection system for in-vehicle can bus communications. *IEEE Access*, 8:185489–185502.
- Khan, Z., Chowdhury, M., Islam, M., Huang, C.-Y., and Rahman, M. (2020). Long short-term memory neural network-based attack detection model for in-vehicle network security. *IEEE Sensors Letters*, 4:1–4.
- Kukkala, V. K., Thiruloga, S. V., and Pasricha, S. (2020). INDRA: Intrusion detection using recurrent autoencoders in automotive embedded systems. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 39:3698–3710.
- Lestyán, S., Ács, G., Biczók, G., and Szalay, Z. (2019). Extracting vehicle sensor signals from CAN logs for driver re-identification. In Mori, P., Furnell, S., and Camp, O., editors, *Proceedings of the 5th International Conference on Information Systems Security and Privacy, ICISPP 2019, Prague, Czech Republic, February 23-25, 2019*, pages 136–145. SciTePress.
- Lin, Y., Chen, C., Xiao, F., Avatefipour, O., Alsubhi, K., and Yunianta, A. (2021). An evolutionary deep learning anomaly detection framework for in-vehicle networks – CAN bus. *IEEE Transactions on Industry Applications*, to appear.
- Lokman, S. F., Othman, A. T., Musa, S., and Husaini, A. B. M. (2019). Deep contractive autoencoder-based anomaly detection for in-vehicle controller area network (CAN). In Husaini, A. B. M., Sidik, M., and Ochsner, A., editors, *Progress in Engineering Technology: Automotive, Energy Generation, Quality Control and Efficiency*, pages 195–205. Springer International Publishing, Cham.

- Mo, X., Chen, P., Wang, J., and Wang, C. (2019). Anomaly detection of vehicle can network based on message content. In Li, J., Liu, Z., and Peng, H., editors, *Security and Privacy in New Computing Environments*, pages 96–104, Cham. Springer International Publishing.
- Moore, M. R., Bridges, R. A., Combs, F. L., Starr, M. S., and Prowell, S. J. (2017). Modeling inter-signal arrival times for accurate detection of can bus signal injection attacks: A data-driven approach to in-vehicle intrusion detection. In *Proceedings of the 12th Annual Conference on Cyber and Information Security Research, CISRC '17*, New York, NY, USA. Association for Computing Machinery.
- Muter, M. and Asaj, N. (2011). Entropy-based anomaly detection for in-vehicle networks. *IEEE Intelligent Vehicles Symposium, Proceedings*, pages 1110–1115.
- Negi, N. S., Jelassi, O., Clemencon, S., and Fischmeister, S. (2019). A LSTM approach to detection of autonomous vehicle hijacking. In *Proceedings of the 5th International Conference on Vehicle Technology and Intelligent Transport Systems - Volume 1: VEHITS*, pages 475–482. INSTICC, SciTePress.
- Novikova, E., Le, V., Yutin, M., Weber, M., and Anderson, C. (2020). Autoencoder anomaly detection on large CAN bus data. In *Proceedings of DLP-KDD 2020*, New York, NY, USA. ACM.
- Remeli, M., Lestyán, S., Ács, G., and Biczók, G. (2019). Automatic driver identification from in-vehicle network logs. In *2019 IEEE Intelligent Transportation Systems Conference, ITSC 2019, Auckland, New Zealand, October 27-30, 2019*, pages 1150–1157. IEEE.
- Song, H. M., Kim, H. R., and Kim, H. K. (2016). Intrusion detection system based on the analysis of time intervals of CAN messages for in-vehicle network. In *2016 International Conference on Information Networking (ICOIN)*, pages 63–68.
- Song, H. M., Woo, J., and Kim, H. K. (2020). In-vehicle network intrusion detection using deep convolutional neural network. *Vehicular Communications*, 21:100198.
- Stone, B. J., Graham, S., Mullins, B., and Kabban, C. S. (2018). Unsupervised time series extraction from controller area network payloads. In *2018 IEEE 88th Vehicular Technology Conference (VTC-Fall)*, pages 1–5.
- Tariq, S., Lee, S., and Woo, S. S. (2020). CANTransfer: Transfer learning based intrusion detection on a controller area network using convolutional lstm network. In *Proceedings of the 35th Annual ACM Symposium on Applied Computing, SAC '20*, page 1048–1055, New York, NY, USA. Association for Computing Machinery.
- Taylor, A., Japkowicz, N., and Leblanc, S. (2015). Frequency-based anomaly detection for the automotive can bus. In *2015 World Congress on Industrial Control Systems Security (WCICSS)*, pages 45–49.
- Taylor, A., Leblanc, S., and Japkowicz, N. (2016). Anomaly detection in automobile control network data with long short-term memory networks. In *2016 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, pages 130–139.
- Tomlinson, A., Bryans, J., Shaikh, S. A., and Kalutarage, H. K. (2018). Detection of automotive can cyber-attacks by identifying packet timing anomalies in time windows. In *2018 48th Annual IEEE/IFIP International Conference on Dependable Systems and Networks Workshops (DSN-W)*, pages 231–238.
- van den Oord, A., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A., and Kavukcuoglu, K. (2016). Wavenet: A generative model for raw audio. *ArXiv*, abs/1609.03499.
- Weber, M., Wolf, G., Sax, E., and Zimmer, B. (2018). On-line detection of anomalies in vehicle signals using replicator neural networks. In *ESCAR 2018*.