



Budapest University of Technology and Economics
Faculty of Electrical Engineering and Informatics
Department of Network Systems and Services

Attack detection on CAN network using TCN

Scientific Students' Association Report

Author:

Beatrix Koltai

Advisors:

András Gazdag
Gergely Ács, Ph.D.

2023

Contents

Kivonat	iii
Abstract	iv
1 Introduction	1
2 Related Work	4
3 Background	6
3.1 CAN	6
3.2 Intrusion Detection Systems	7
3.3 Temporal Convolutional Networks	7
4 Attacker model	9
5 Proposed solution	10
5.1 Preprocessing of CAN Traffic	10
5.2 Grouping of Correlated Signals	10
5.2.1 Correlation analysis	10
5.2.2 Clustering of signals	12
5.3 Signal Forecasting	14
5.4 Decision	15
5.4.1 Determining the parameter of ℓ	15
5.5 Discussion	17
5.5.1 Why Grouping Correlated Signals?	17
5.5.2 Cost Analysis, Scalability and Robustness	17
6 Evaluation	19
6.1 Dataset	19
6.2 Model Architecture and Parameters	20

6.3	Comparison with Baselines	21
6.4	Evaluation Metrics	22
6.5	Results	23
7	Conclusion	26
	Acknowledgements	27
	Bibliography	28
	Appendix	31

Kivonat

Az elmúlt években jelentősen megnőtt a járművekben található elektronikus vezérlők száma. Ez a változás, a külső interfészek elterjedésével reális veszéllyé tette a kibertámadásokat. Mivel egy támadó akár már távolról is hozzáférhet a járműhöz, annak belső rendszerének a sérülékenységeit kihasználva akár át is veheti az irányítást a jármű felett. Ennek a problémának egyik oka, hogy a CAN protokoll (a legszélesebb körben használt protokoll járművek belső hálózatában) nem tartalmaz semmilyen biztonsági mechanizmust, és így egy támadó potenciálisan megtámadhatja az elektronikus vezérlőket ezen keresztül. A legsúlyosabb problémát a CAN üzeneteket módosító támadások jelentik, amelyek detekciója több szempontból is kihívást jelent. A támadások észlelésére szolgáló rendszerek (Intrusion Detection Systems, IDS) alkalmazása egy lehetséges megoldás a járműhálózatokat érintő veszélyek felderítésére és elhárítására. A korábbi megoldások képesek a járművek jeleinek nagyobb rendellenességeit észlelni, viszont a kisebb eltérést okozó támadások detektálásához kifinomultabb megoldásokra van szükség. Egy olyan anomália-felismerő rendszert mutatunk be, amely a jelek normális viselkedésére és az azok közötti kapcsolatokra (korrelációra) együttesen épít.

A fő probléma, amellyel ez a dolgozat foglalkozik, hogy a korábban javasolt detekciós eljárások hiányosságot mutatnak két konkrét támadási helyzet esetében. Egyrésztől, amennyiben egy támadó olyan módosítást hajt végre, amely más helyzetben normális viselkedés is lehetne, akkor azt kizárólag viselkedés alapú módszerrel nem lehet jól felismerni. Másrésztől, amennyiben a támadó ismeri a jelek korrelációját és képes a jeleket egyszerre módosítani, akkor ezt megteheti úgy, hogy azok korrelációja továbbra is megmaradjon. Ezesetben a pusztán korrelációra támaszkodó megoldások nem fogják detektálni a támadást.

Egy olyan megközelítést javasolunk az üzenetmódosítások észlelésére, mely a két módszert kombinálja: az idősor-előrejelzés és a jelkorreláció elemzésének együttes használatával elemzi a CAN forgalmat. A két megközelítés kombinálásával megoldjuk az egyes módszerek egyedüli alkalmazása során fellépő hiányosságokat.

Az adatok időbeli függőségeinek modellezésére egy több csatornás Temporal Convolutional Network idősoros hálót használunk, melyet a jelek egy csoportjára tanítunk be. A csoportosítás a jelek korrelációja alapján történik, így ezen jelek együttes predikciója során azok korrelációját is figyelembe vesszük. A csoporton belüli jelekhez mind előállítunk egy predikciót, mely így függeni fog a jel viselkedésétől, és a csoporton belüli korrelációtól, tehát bármelyik jellemzőben történne az anomália, a detekció során lesz róla információnk.

A detekció során olyan összehasonlítást alkalmazunk, mely figyelembe veszi, hogy milyen mértékben tért el az előrejelzés a tényleges értéktől, és azt, hogy ez az állapot milyen hosszán állt fenn.

A dolgozatban bemutatjuk a módszer hatékonyságát egy már meglévő, korábbi megoldások tesztelésére is használt CAN forgalmi támadásokat tartalmazó adathalmazon.

Abstract

The number of electronic controls in vehicles has increased significantly in recent years. This change, combined with the widespread of external interfaces, has made cyber-attacks a real threat. As an attacker can now gain access to a vehicle remotely, they can exploit vulnerabilities in the vehicle's internal systems to take control of the vehicle. One of the reasons for this problem is that the CAN protocol (the most widely used protocol in the vehicles' internal networks) does not contain any security mechanism. Through this, an attacker could attack electronic control units. The most severe problem is CAN message modification attacks, which are challenging to detect in several ways. Intrusion Detection Systems (IDS) are a possible solution to detect and mitigate threats to vehicle networks. Previous solutions can detect significant anomalies in vehicle signals, but more sophisticated solutions are needed to detect attacks that cause only minor deviations. We present an anomaly detection system that relies on the expected behavior of the signals and the relationships (correlation) between them.

The main problem addressed in this paper is that the previously proposed detection methods indicate shortcomings in two specific attack situations. On the one hand, if an attacker makes a modification that could be a normal behavior in another situation, it cannot be well detected by a behavior-based method alone. On the other hand, if the attacker knows the correlation of the signals and can modify all the signals simultaneously, they can do so while maintaining their correlation. In this case, detection methods that rely on correlation alone will not detect the attack.

We propose an approach to message modification detection that combines the two methods: time series prediction and signal correlation analysis. Combining the two approaches, we overcome the shortcomings of using each method alone.

To model the temporal dependencies of the data, we use a multichannel Temporal Convolutional Network trained on a specified group of signals. The clustering is based on the signals' correlation, so the joint prediction of these signals will correspond to their correlation. For each signal within the group, we produce a prediction, which will then depend on the behavior of the signal and the correlation within the group so that if an anomaly occurs in any of these features, we will have enough information to detect it.

The detection process uses a comparison that considers how much the prediction deviated from the actual value and the length of time this state lasted.

This work demonstrates the method's effectiveness on an existing dataset of CAN traffic attacks, which has been used to test previous solutions.

Chapter 1

Introduction

Securing vehicular communication networks is becoming crucial as the automotive industry rapidly evolves and increasingly adopts connectivity. Applying Intrusion Detection Systems (IDS) in specific domains is becoming essential for identifying and mitigating threats to vehicular networks. One such domain is the vehicles' inner communication on the Controller Area Network (CAN).

The CAN bus is a complex network of Electronic Control Units (ECUs) that collaborate to provide the necessary functions of the vehicle. Cyber attacks targeting these ECUs can have dire consequences for safety-critical subsystems such as brakes, the engine, or the steering wheel. A malfunctioning vehicle not only endangers passengers and others around it but also impacts the VANET (Vehicular Ad-hoc Network). Compromising data used in Vehicle-to-Everything (V2X) communication, an attacker could spread malicious information and alter the behavior of others, which could cause congestion or severe accidents in an urban environment. An attacker can have financial motivation besides deteriorating reliability and driving safety. Gaining control over the vehicle could allow theft, stealing sensitive data, and sabotaging the system.

Since the CAN protocol does not implement any security measures [3], an attacker can potentially attack the ECUs by making communication inaccessible, injecting new malicious messages, or even modifying valid messages. DoS (Denial-of-Service) attacks disable the benign CAN communication by flooding the network with the highest priority messages. However, this attack can be easily detected because the network load is significantly increased during the attack. Message injection can also affect specific vehicle functions, but these attacks are also easy to detect, with simple statistical methods, as injected messages cause a recognizable change in the regular arrival times.

The most challenging issue is message modification attacks that do not introduce new messages to the network, only the data contents are changed. This attack is the hardest to detect due to the variability in traffic patterns, lack of authentication or encryption, the existence of stealthy attack techniques, and the lack of attack signatures. In general, only the continuously changing message data can be used for identifying anomalies that requires general, accurate methods to differentiate between normal and malicious behavior.

After extracting signals from the messages, the detection of malicious message modifications follow two main approaches: time-series forecasting [16], [17], [6] and signal correlation analysis [14], [21]. In time-series forecasting, a machine learning model is trained per signal that predicts the next, expected signal value. Anomaly is reported when there is a substantial deviation between the prediction and the actual value. Unfortunately, this method is incapable of identifying modifications that fall within the usual, non-anomalous range of signal values, even if they constitute an attack. For instance, this limitation is evident when the speed value is modified, causing it to marginally fall below the speed

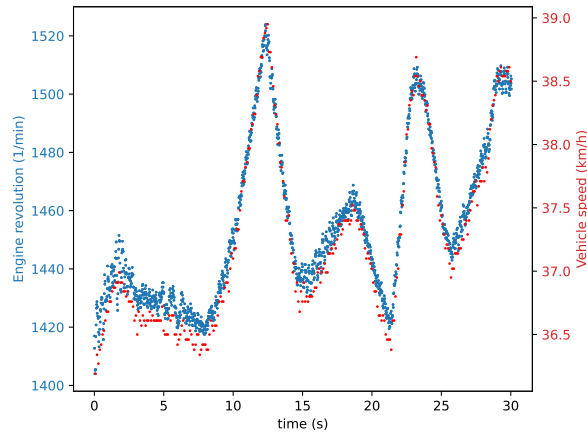


Figure 1.1: Example benign CAN signal (S-1-4).

limit. To overcome this shortcoming, the deviation of the correlation between each pair of signals is checked, where correlation is calculated based on the most recent few minutes' worth of signal data [14], [21]. Indeed, increasing the speed should naturally result in a corresponding increase in the RPM signal; otherwise their correlation would appear anomalous, as Figure 1.1 shows. Consequently, to evade detection, an attacker would need to maintain the original correlation intact and simultaneously modify all correlated signals, which could be prohibitively expensive in practice. Nonetheless, unlike time-series forecasting, this purely correlation-driven approach is unable to identify malicious alterations in signals that lack any correlation between them.

Our proposal combines the merits of both time-series forecasting and correlation analysis, as shown in Figure 1.2. We simultaneously forecast multiple correlated signals and flag an anomaly if the cumulative difference between the predicted values and the actual values of all correlated signals exceeds a specified threshold. The underlying idea is that, as a single model forecasts multiple highly correlated signals, any alteration in one signal will inevitably influence the predictions of all other correlated signals. In other words, we leverage signal correlation not only for more accurate prediction, but also to induce detectable deviation of the predicted signals from the actual ones even if only one of them is maliciously modified. For example, the larger the speed the larger the RPM value, which means that increased speed with constant RPM is likely to produce a noticeable cumulative prediction loss over *both* signals if they are predicted jointly by a single model. Furthermore, unlike pure correlation-based approaches, our method is capable of identifying malicious alterations in signals, even those that lack correlation, when their predicted values deviate significantly from their actual values. Additionally, it can detect attacks in which the attacker modifies correlated signals simultaneously without altering their correlation, yet still induces abnormal behavior.

Our contributions in this work are as follows:

- We employ a combination of time-series forecasting and signal correlation analysis to identify anomalies in the vehicular CAN bus. Our unsupervised method relies solely on *unlabeled* CAN traces for training and calibration prior to deployment. It operates by simultaneously predicting correlated signals that allows a more accurate detection of abnormal behaviour.
- We assess the effectiveness of our approach using a dataset comprising eight distinct message modification attack types. Our results demonstrate a substantial performance improvement over the state-of-the-art: we achieve a detection rate of 95% (compared to 68%) with a precision of 80% (versus 30%). Additionally, our method exhibits a minimal average detection delay of just

0.38 seconds.

- Finally, we show that in addition to modification attacks, our solution also effectively identifies injection attacks, allowing the identification of both types of attacks by a single algorithm.

The rest of the paper is organized as follows: Chapter 2 briefly covers prior research and developments in anomaly detection in Controller Area Networks. Chapter 3 summarizes the relevant background of the CAN bus and vehicular intrusion detection solutions. The attacker model is introduced in Chapter 4. Chapter 5 describes the proposed anomaly detection mechanism, the training process, and the detection process. Chapter 6 evaluates the performance of the method on real-world CAN data. Finally, in Chapter 7 we conclude our paper.

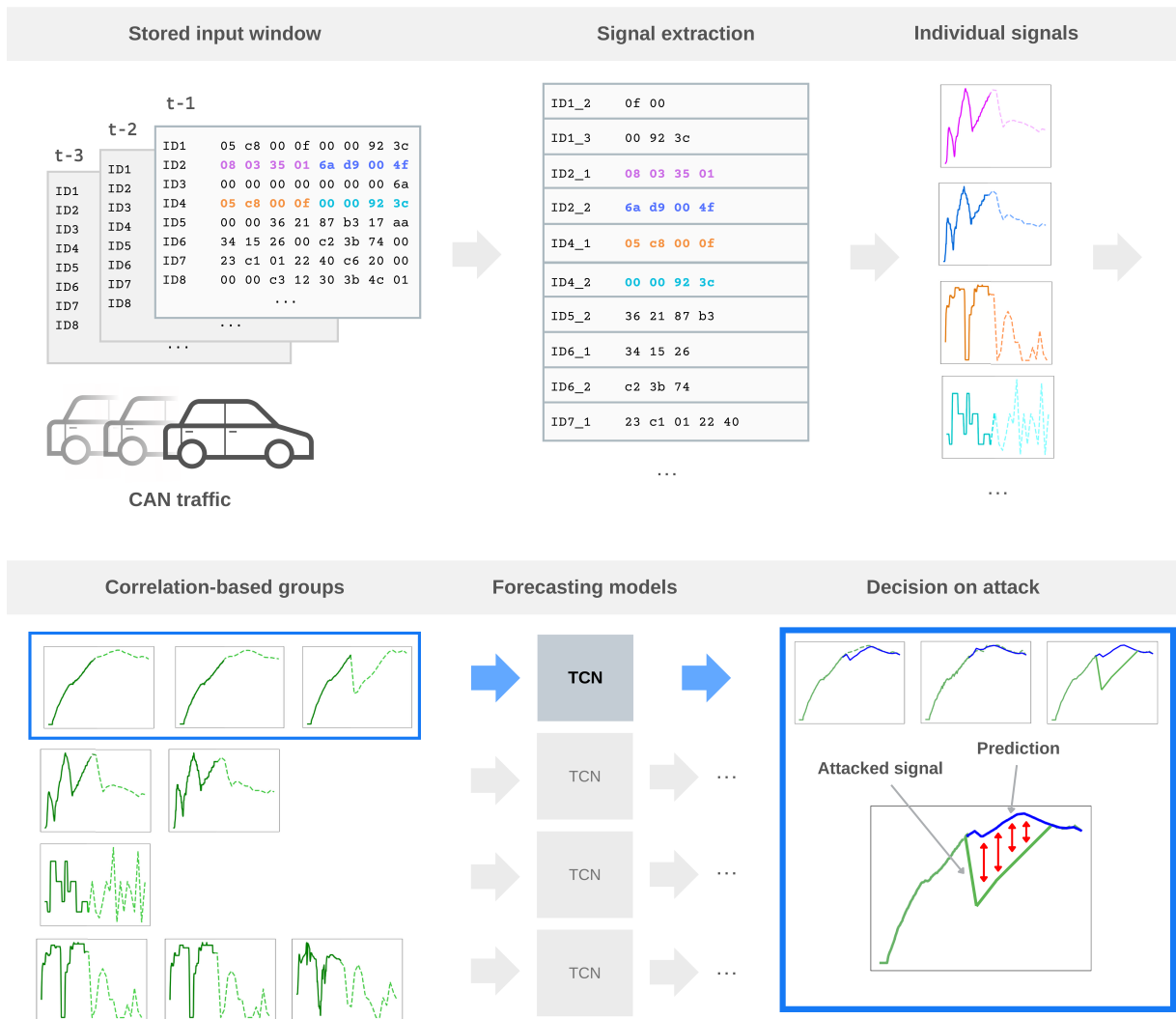


Figure 1.2: High-level layout of our correlation-based approach.

Chapter 2

Related Work

Intrusion detection systems used in in-vehicle networks differ from those used on the Internet because there are limited known attack signatures. Most research results are based on unsupervised learning, as the available data can only be used appropriately to describe the benign state of the systems. Following this approach, papers have been published on detecting message injection and modification attacks.

IDS systems often rely on measuring and monitoring the timestamp of message arrivals to detect injection attacks. Due to the periodical timing of CAN data messages in a benign state, timing-based detection methods can effectively detect message insertions and drops [26, 13]. Young et al. showed that the constant nature of the inter-arrival times can also change for short periods of time during transitions of vehicle state [28]. They propose analyzing the message arrival times in the frequency domain to build a robust detection algorithm even for state transitions. In their research, Müter et al. proposed measuring the message entropy for anomaly detection [22]. While this approach successfully detected injection attacks, they also demonstrated the shortcomings of their approach in short-duration attack scenarios. Machine learning has also been used for the detection of injection attacks. Guidry et al. have proposed using a one-class classification method [15]. Features of their model included inter-arrival times, the transmission frequencies, and the deviations from the typical inter-frame times. They measured the effectiveness of different one-class classification-based approaches and concluded that the S-SVDD method performs the best with an average of 85% detection rate.

Attackers, however, cannot only inject messages into the bus, but it is also possible for them to modify messages, as described in Chapter 4.

In [18], the proposed method can detect these modification attacks by utilizing the transient state at the beginning of a modification attack. For a short time missing messages could indicate a suspension attack as a preparation step for a modification attack. However, if this phase is not detected in time, the rest of the attack will be successful.

In recent years, many papers have been published on identifying modification attacks based only on the message data contents. Among others, researchers tackled the problem by continuously measuring the relationship between data fields, forecasting future data values and later identifying deviations between the predictions and actual values.

CAN signal correlation analysis is proposed in [14] to identify modification attacks. Even though this approach is robust against attacks that target highly correlated signals, its effectiveness is generally limited. The proposed solution calculates correlations between signals regularly in two different time windows to identify ongoing anomalies. In [21], the authors extend correlation analysis with hierarchical clustering. Their results are demonstrated on a dataset, but it is not compared to other baseline results. As the

presented framework can only handle entire traffic logs, it is not applicable as a real-time detector for the CAN bus but only as a forensics tool.

Time series forecasting is also used to predict future values in CAN communication, either on message or signal level. These predictive methods can identify possible modification attacks by measuring deviations between predicted and actual measured values.

Using a neural network for anomaly detection has been proposed in CANet [16]. The authors used independent LSTM models for each message ID to capture the corresponding signal's temporal dynamics and forecast its future values. The output of all models is then fed into a fully connected autoencoder layer, allowing the network to consider the interdependencies of signals. Although this approach exploits relations between signals for detection, this information is not directly used in the network structure. In [17], the INDRA framework was proposed, which analyzes temporal patterns and behavior of messages using Gated Recurrent Unit (GRU) based recurrent autoencoders. One such autoencoder was trained for each message ID to reconstruct signals within the message. The authors show that INDRA outperforms CANet in accuracy and false positive rate. In [6], the authors introduce a Temporal Convolutional Network based detection system. Their approach separates CAN signals and builds individual predictor models for each signal, similar to CANet and INDRA. However, as TCN networks are smaller and faster than previous neural networks, such as LSTMs, their solution outperforms all previous results. In this paper, we improve on the TCN-based approach by introducing signal clustering to improve detection results while reducing the mechanism's footprint.

Chapter 3

Background

This section provides an overview of the CAN network’s operation within vehicles, outlines the typical methods used to build an Intrusion Detection System, and introduces the application of Temporal Convolutional Neural Networks (TCNs) along with signal correlation analysis as part of our proposed anomaly detection approach.

3.1 CAN

Modern-day vehicles have a complex internal control system comprised of ECUs, each assigned to manage a specific function. These ECUs are interconnected via networks, the most important being the Controller Area Network. While this system has proven reliable over the years, external interfaces have exposed it to potential attacks [5, 1].

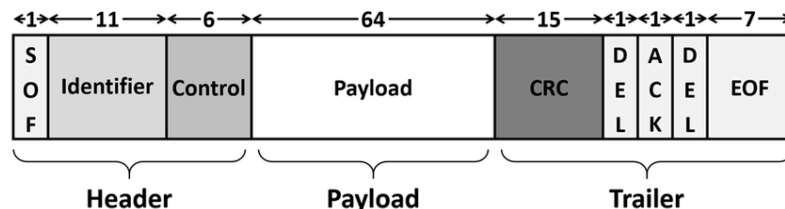


Figure 3.1: Structure of a CAN frame [17].

On the CAN bus information is transmitted in frames. A CAN frame is shown in Figure 3.1, containing header, payload, and trailer segments. The header contains the start of the frame signal for synchronization, the message identifier (ID), and the data length code (DLC), which specifies the payload’s length. The actual data to be transmitted is in the payload segment. The trailer segment is mainly used for error checking at the receiver’s end. The cyclic redundancy check (CRC) is used for the data integrity check, while the acknowledgment (ACK) is used to confirm reception.

Messages sent over the CAN network have an ID, either 11 bits or 29 bits long. A typical passenger vehicle uses an 11 bit identifier. The data section can range from 0 to 8 bytes of data. Within the data part, various digital and analog signals are encoded. Manufacturers do not disclose how the signals are encoded, but they can be reverse-engineered using methods previously proposed in the literature [19, 27, 20].

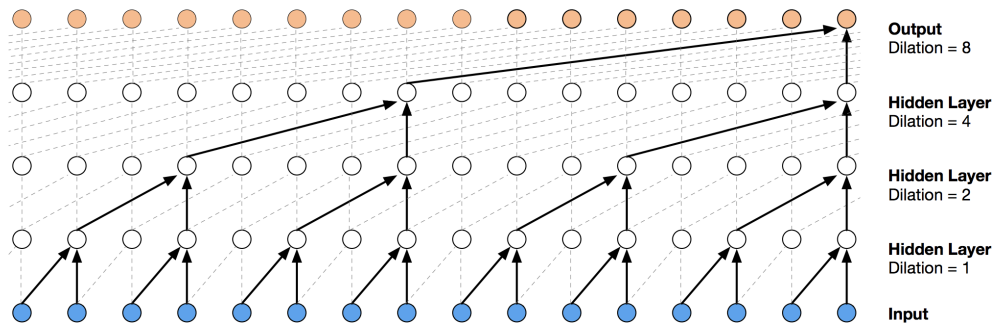


Figure 3.2: Structure of a stack of dilated causal convolutional layers in TCN [24].

3.2 Intrusion Detection Systems

In order to detect attacks, Intrusion Detection Systems (IDS) mainly utilize two methods: signature-based and anomaly-based detection [2].

Signature-based detection systems search for specific attack features in the examined traffic. While they have a low false positive rate, they require knowledge of the attacks to detect them accurately. Any attacks that are not modeled in the signature database will not be detected by the system.

An anomaly detection system relies on learning the system's normal behavior and identifying any messages that indicate a deviation from this benign state. This approach is beneficial in situation where it is not possible to describe the attacks in advance.

Vehicular networks show large variations, as manufacturers significantly change the built-in features between vehicle types. In this complex landscape, attacks are also customized for each target. Thus, creating a comprehensive database of every attack is not feasible, therefore vehicular attack detection systems are mostly anomaly-based.

Our detection model is based on an unlabeled data model, built from benign network traffic of a test vehicle, that implements an anomaly-based IDS system. Although the dataset we use includes real attacked CAN data, it will only be used for testing and evaluation purposes as it is not representative of all possible attack types.

3.3 Temporal Convolutional Networks

Convolutional Neural Networks (CNNs) and Temporal Convolutional Networks (TCNs) are deep learning architectures widely used for various tasks, including image recognition and natural language processing. They offer significant benefits when applied to time series data, making them suitable for detecting anomalies in the Controller Area Network (CAN) [6].

CNNs are designed to process grid-like data, such as images, by applying convolutional filters to extract spatial features. In the case of time series data, 1-dimensional causal convolutions can be used to identify local patterns and dependencies within the data.

A TCN is a type of deep learning architecture designed explicitly for sequential data, such as time series. To process sequences in parallel, TCNs use dilated convolutions, which enable them to capture long-

range dependencies efficiently, as shown in Figure 3.2. This ability is critical in identifying anomalies that may occur over extended periods or exhibit complex temporal behaviors. Additionally, TCNs stack multiple layers for hierarchical feature extraction. They also employ causal padding, ensuring only past and present information is used. Due to these features, TCNs are suitable for various applications, including time series forecasting and anomaly detection.

TCNs can handle large volumes of data, making them suitable for analyzing extensive CAN message traffic. This architecture can be optimized for real-time processing, allowing immediate anomaly detection and response in safety-critical CAN systems.

Chapter 4

Attacker model

This section discusses the attacker model and the attack surface of a CAN network. We describe the capabilities and goals of an attacker and classify the potential attacks that an attacker may perform on CAN messages.

We assume that the attacker can gain access to the vehicle using the most common attack vectors [5]. The goal of the attacker is to send forged data to an ECU, forcing it into a corrupt state. This could cause problems anywhere between showing invalid values on the dashboard to making the vehicle completely unusable or stealing it¹, depending on the target ECU. This goal can be achieved in multiple ways. An attacker with physical access to the vehicle can add new devices to the CAN network. Vehicles with wireless interfaces, such as Bluetooth, WiFi, or a 3G/4G/5G connection, can also be attacked remotely. After exploiting a vulnerability in the communicating ECU, similar CAN transmission capabilities can be gained. This is the first necessary step of any attack against the CAN bus. The CAN network operates reliably under normal conditions; however, due to the absence of security provisions within its specification, it remains susceptible to potential attacks. Once an attacker has the capability to interact with the CAN bus, there are multiple possible attack strategies, including DoS, message injection, and message modification. The latter two are also referred to as a fabrication and a masquerade attack.

First, we focus on the most challenging problem, which is the message modification attack. During these attacks the repetition times of the messages are unchanged, as there are no new messages introduced to the network. Hence, messages arrive at their expected time but with a modified data content. Carrying out such an attack requires strong technical skills, nevertheless, its feasibility has already been demonstrated in [7]. A practical implementation of such an attack exploits the error handling mechanism of the CAN protocol. If a device detects an error during transmission, an error signal bit can be used to inform the sender about the problem. Repeated error signals can force an ECU into an error state. In this state all further message transmissions are suspended, allowing an attacker to take the place of the ECU in the communication and send modified messages. Therefore, identifying modification attacks based only on meta-data (e.g., the number or timing of CAN messages) is not possible. In this paper, we present a novel anomaly detection mechanism, designed to detect such attacks.

We also evaluate the performance of our model for message injection attacks. A less sophisticated detection mechanism, based only on statistical properties of the message arrival times, can also detect such attacks. However, if the modification attacks require a neural network-based anomaly detector, it is beneficial to know the detector's performance for every scenario.

¹<https://arstechnica.com/information-technology/2023/04/crooks-are-stealing-cars-using-previously-unknown-keyless-can-injection-attacks>

Chapter 5

Proposed solution

Our solution has three main components: after extracting signals from the raw CAN traffic, (1) correlated signals are grouped together using clustering, (2) a separate and independent supervised forecasting model per group predicts the next value of all correlated signals within a group, and finally (3) an anomaly is reported if at least one of the forecasting model's predictions deviate significantly from the true, observed values of the predicated signals. We detail the operation of each component as follows.

5.1 Preprocessing of CAN Traffic

All signals from the available CAN messages are extracted using the manufacturer's specification or any state-of-the-art automatic extraction tool [23, 19, 27]. As not all extracted signals are equally useful for anomaly detection, a subset K of all extracted signals are retained while the rest are dropped. Indeed, useless signals are extracted from unused parts of the CAN messages (i.e., there is no device in the vehicle that uses that part of the message), or carry constant values with no predictive power. This filtering process also helps minimize the size of the forecasting model detailed in Section 5.3. Finally, all retained signals are normalized by dividing each signal value by their theoretical maximum that is either specified by the manufacturer, or computed as $\lceil 2^s \rceil$ where s is the number bits used to store the signal in the CAN message.

5.2 Grouping of Correlated Signals

All retained K signals are clustered into C groups based on their pairwise correlation value. Specifically, each signal is first assigned to a separate cluster and then the closest clusters are iteratively merged until the number of clusters attains K , where the closeness of two clusters is measured by a chosen correlation metric of their respective centroids. Our approach is not restricted to any specific similarity measure or clustering technique. Still, as we show in this section, and in Chapter 6, linear correlation with hierarchical clustering is effective in practice.

5.2.1 Correlation analysis

We have analyzed different correlation metrics to support the claim that linear correlation is a viable option for clustering signals. Our investigation included three correlation methods, namely Pearson,

Spearman, and Kendall correlation, which are commonly used for evaluating the strength of relationships between two variables.

Pearson’s correlation coefficient [9] measures the linear relationship between two continuous variables, suitable for typical analog signals on the CAN bus, such as speed, PRM, etc. The formula for calculating r_p Pearson’s correlation can be seen in Equation 5.1, where X and Y are the two variables, and $x_j, y_j, j = 1, 2, \dots, n$ the observed data points. Pearson correlation assumes that the variables follow a normal distribution and have a linear relationship. Correlation values range from -1 to 1, where -1 indicates a perfect negative linear relationship, 1 indicates a perfect positive linear association, and 0 shows no linear relationship. It is sensitive to outliers, meaning extreme values can significantly influence the correlation value.

$$r_p = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (5.1)$$

Spearman Rank Correlation [9] measures the strength and direction of the monotonic relationship between two variables by calculating the correlation based on the ranks of data points. It is calculated in the same manner as the Pearson correlation, except that the Spearman correlation is calculated after the two variables have been ranked and transformed to values between 1 and the number of observations. Other than a monotonic relationship, it does not assume linearity or follow any specific distribution. The ranking property of this metric makes it robust to outliers. Similarly to the Pearson correlation, resulting values range from -1 to 1, with -1 indicating a perfect inverse rank correlation, 1 indicating a perfect rank correlation, and 0 indicating no rank correlation.

Like Spearman, Kendall’s Tau, also known as Kendall’s rank correlation coefficient [9], does not assume linearity or follow any specific distribution. Values range from -1 to 1, where -1 indicates a perfect inverse ordinal association, 1 indicates a perfect ordinal association, and 0 indicates no ordinal association. Kendall’s Tau is often considered more robust than Spearman’s when dealing with tied data values.

A heatmap for each correlation method displaying the pair-wise correlation between signals is shown in Figure 5.1. The magnitude of the correlation is represented by the darkness of the color and the size of each point. The sign of the correlation is encoded in the hue of the color.

As can be seen from Figure 5.1, all three correlation methods identify almost the same correlation. However, some significant differences occur in the case of signals 0290_1, 0410_4, and 300_4, which correlate only with signals 0290_4 and 0290_2 when calculating the Pearson method. Furthermore, there is a stronger correlation between the group of signal 0120_0 and signal 0120_1 and the group of signal 0110_1 and signal 0110_3, as observed by the Pearson method, whereas a weaker correlation can be measured for Spearman and Kendall. We manually reviewed these discrepancies and determined no real connections between these signals.

Our measurement results indicate that the Pearson correlation method is the most suitable for our needs. However, it is worth noting that the final grouping of signals has a significant impact on the detection accuracy rather than the correlation results directly. As such, we are evaluating all three correlation computation methods while exploring clustering algorithms to determine the optimal signal grouping for our purposes.

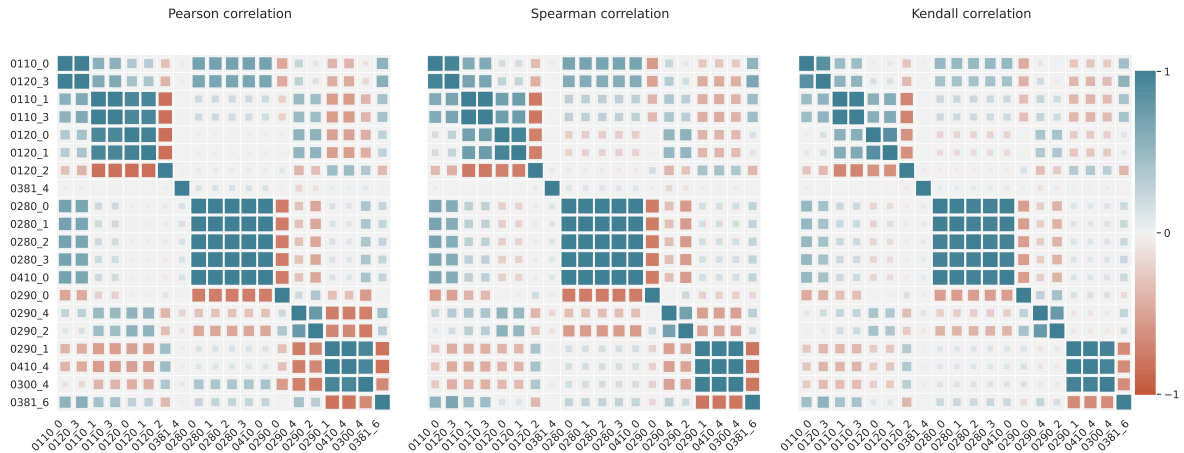


Figure 5.1: Heatmaps of different correlation metrics used to determine similarity between signals (Pearson, Spearman and Kendall correlation form left to right). To better illustrate the magnitude of the correlation, we also varied the size of each point on the heatmaps, which is proportional to the darkness of the color. The sign of the correlation is encoded in the hue of the color.

5.2.2 Clustering of signals

We use the measured correlations to group the signals in the prediction phase. This clustering has two benefits. First, it is more efficient to treat similar signals together, as we need fewer models compared to predicting each signal alone. Second, it is more accurate due to the additional information about interdependencies. The joint prediction of the groups will thus implicitly exploit the correlation of the signals since if it changes, the group’s forecast will also change.

As previously mentioned, our approach is not limited to any particular clustering technique. However, to demonstrate the effectiveness of our chosen hierarchical clustering, we compared it to four other clustering techniques, each tested with all the correlations discussed previously.

We compared four distinct clustering algorithms on our dataset - DBSCAN, Affinity Propagation, Hierarchical Clustering, and Mean Shift Clustering. We chose only clustering techniques that do not require the number of clusters to be specified in advance, as we do not know how many groups can be utilized.

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a density-based clustering algorithm that groups data points based on their density within the dataset [25]. It can discover clusters of arbitrary shapes and is robust to outliers called noise points. The method identifies core points (dense data points), border points (points on the edge of a cluster), and noise points (isolated points).

Affinity Propagation is an exemplar-based clustering algorithm that selects a set of data points as exemplars and assigns the rest of the points to the nearest exemplar [11]. It automatically determines the number of clusters by choosing exemplars in the data. Affinity Propagation can be sensitive to the choice of similarity or distance metric, and the number of exemplars can significantly affect the results. We observed this when we computed the Affinity Propagation method using Pearson, Spearman, and Kendall metrics, as it produced quite different clusters.



Figure 5.2: Comparing signal groups produced by different clustering techniques, with different correlation metrics. From top to bottom, the result for 1) MeanShift clustering with Spearman correlation, 2) Hierarchical clustering with Pearson correlation, 3) DBSCAN clustering with Kendall correlation and finally 4) Affinity propagation clustering with Spearman correlation can be seen. First column indicates the name of the chosen pair of clustering and correlation method, second column shows a 2D representation of the clustering, while the first column shows each signal, with the same background color for each signal in a group.

Hierarchical clustering builds a tree-like hierarchy of clusters, often represented as a dendrogram [25]. It can be agglomerative (bottom-up) or divisive (top-down). Agglomerative clustering starts with individual data points as clusters. It merges them iteratively, while divisive clustering begins with a single set containing all data points and splits them into smaller groups [10].

Mean Shift clustering is a mode-seeking clustering algorithm that aims to find the modes or peaks of data density [4]. It is beneficial for finding clusters with non-uniform shapes or densities. Mean Shift is sensitive to the bandwidth parameter, which affects the size and shape of the groups.

To visually represent signals and their relationships in a 2D figure, we use dimensionality reduction. This involves projecting the signals into a 2D space while keeping their similarity or dissimilarity. To achieve this, we use Multidimensional Scaling (MDS), which aims to position data points in a way that preserves their pairwise distances or similarities as accurately as possible in the reduced-dimensional space. Each point in the 2D subfigure corresponds to a signal.

This representation is purely for visualization purposes. It does not impact the chosen correlation method or clustering technique.

We compared the aforementioned clustering techniques by analyzing the resulting groups using all three correlation metrics. The result of the executions can be seen in Figure A.1. Here, we evaluate in details the clustering results with Pearson correlation, which is shown in Figure 5.2.

During the DBSCAN execution, we had to adjust most of the algorithm parameters. It was evident that their values greatly affected the output. We noticed that some signals were assigned to separate groups, even though they seemed to belong together. Due to this issue, we decided to exclude this method from further consideration.

We also assessed the Affinity Propagation method but found it too sensitive to parameters. Even after fine-tuning the parameters to obtain the best possible results, we found that it was still grouping signals that did not belong together. Therefore, we decided to exclude this method as well.

We then evaluated MeanShift and Hierarchical clustering algorithms and found that both methods essentially gave us the same results. After manually reviewing the output, we determined that the outcome aligned with our expectations. As such, we concluded that both methods were suitable for our use case. Ultimately, we opted to use the Hierarchical clustering algorithm for ease of use.

5.3 Signal Forecasting

We train C supervised models on the clustered CAN data in order to predict the next upcoming signal value: all retained K signals are divided into equally-sized overlapping segments using a sliding window with size w , and each segment serves as input to the forecasting model to predict the subsequent signal value immediately following the segment.

More precisely, let a signal with ID s be represented as a time series (T_1^s, \dots, T_n^s) after pre-processing, and $\mathbf{M}^G = [(T_1^{g_1}, T_2^{g_1}, \dots, T_n^{g_1}), \dots, (T_1^{g_{|G|}}, T_2^{g_{|G|}}, \dots, T_n^{g_{|G|}})] \in \mathbb{R}^{|G| \times n}$ denotes the time series of all correlated signals in group G , where $G = \{g_1, \dots, g_{|G|}\}$ are the set of signal IDs belonging to G . For any signal group G , a forecasting model f_G simultaneously predicts the next element of each signal of the group: given the most recent w signal values $\mathbf{M}_{t-w:t}^G = [(T_{t-w}^{g_1}, T_{t-w+1}^{g_1}, \dots, T_{t-1}^{g_1}), \dots, (T_{t-w}^{g_{|G|}}, T_{t-w+1}^{g_{|G|}}, \dots, T_{t-1}^{g_{|G|}})] \in \mathbb{R}^{|G| \times w}$ as input, the forecasting model predicts the next value $\mathbf{M}_{t:t+1}^G = (T_t^{g_1}, T_t^{g_2}, \dots, T_t^{g_{|G|}})^\top \in \mathbb{R}^{|G|}$ of every signal in G . Before deployment, all forecasting models are trained on CAN data that comes from the same or sufficiently similar distribution as the actual CAN traffic after deployment.

5.4 Decision

We compare the prediction made by every forecasting model with the actual, observed values of the signals, and report anomaly if the deviation of the prediction is too large for any group.

More precisely, let $\mathbf{O}_{t:t+1}^G$ denote the actual, observed value of the signals at time t in group G after performing the pre-processing steps detailed in Section 5.1. The prediction error for group G at time t is defined as

$$\text{err}_G(t) = \frac{1}{|G|} \|f_G(\mathbf{O}_{t-w:t}^G) - \mathbf{O}_{t:t+1}^G\|_2^2 \quad (5.2)$$

which measures the mean squared error (MSE) between the actual signal values and the values predicted by f_G from the last w observed values of the signal. Note that \mathbf{O} denotes the true value of the signal that is observed on-line after the deployment of the trained forecasting model f_G .

A naive method of detection is to directly compare the prediction error with a threshold τ , and report anomaly if $\text{err}_G(t) \geq \tau$ for any group G . However, since the variance of $\text{err}_G(t)$ can be large depending on the accuracy of the forecasting model f_G , this approach can yield large detection error: any value of τ would induce either too many false positives (for smaller τ) or false negatives (for larger τ). To mitigate such effect of forecasting inaccuracy, we rather compare the mean of the last ℓ error values with the threshold, that is, report anomaly if $(1/\ell) \sum_{i=t-\ell}^{t-1} \text{err}_G(i) \geq \tau$ for any group G . This approach also more reliably detects stealthier attacks that span multiple time slots and involve insignificant modification of the signal value per slot, but surpass the threshold when aggregated.

To adjust τ , we follow the standard three-sigma rule and set τ to three times the standard deviation of $(1/\ell) \sum_{i=t-\ell}^{t-1} \text{err}_G(i)$ plus its expected value on normal (attack-free) traffic [8]. The underlying assumption is that, without adversarial manipulation, the cumulative prediction error lies within three standard deviations of its mean that has a probability of 0.9973 if it is normally distributed (which is the case if ℓ is sufficiently large). The three-sigma rule is applicable even without access to attacked traffic before deployment, otherwise an optimal calibration of τ follows from the Neyman-Pearson lemma.

5.4.1 Determining the parameter of ℓ

We performed experiments to determine the value of ℓ , representing the number of the last loss values considered in the decision-making process. We conducted this experiment by analyzing the detection results for different values of ℓ . A reasonable range for the ℓ values is from 1 to 500, where $\ell = 1$ means we do not calculate mean values, and 500 corresponds to roughly 0.8 milliseconds of data. The goal is to calculate the mean of the loss values in a window that is only on the order of the noise and not yet comparable to a potential attack.

For each ℓ value between 1 and 500, we replaced each loss value with the average of the last ℓ loss values. The distribution of the resulting values is shown in Figure 5.3 on the left side for an example REPLAY attack. The figure shows the effect of two different ℓ values. The first row depicts $\ell = 1$, meaning no averaging was used, and the second row illustrates the case of $\ell = 200$.

The distribution plots on the left show the distribution of loss values containing the attack in blue and the distribution of benign loss values in pale orange. The difference for each bin is highlighted in darker green. It can be seen that at the location of small loss values, there is a lot of error, which is considered noise. Additionally, at larger loss values, there appears to be a more significant error, which is caused by the attack. In the figure, a vertical dashed red line indicates the threshold calculated for the actual loss values. It is important to note that in the case of the blue distribution, the error to the left of the threshold

is not a false negative, which would mean an undetected attack but the same noise as in the benign case.

The plots on the right in the figure show the detection for the given value of ℓ . The vertical orange lines indicate the locations where detection occurred, while the single green bar in the middle represents the location of the actual attack. It can be seen that as ℓ increases, false positives decrease, so there will be fewer detections outside of the attack.

We computed accuracy, false positive rate, and attack detection delay for different values of ℓ . Figure 5.4 illustrates that the false positive rate decreased immediately, reaching a minimum at approximately $\ell = 50$, and then started to increase again. The accuracy increased proportionally and then began to decline around $\ell = 100$.

To ensure that the characteristics of available attacks do not influence the detection process, the value of ℓ was adjusted to minimize false positives in benign signals while avoiding the potential range of an attack.

As discussed in this section, we applied a threshold to the difference between predicted and observed values in our modeling. The choice of threshold and the size of ℓ would depend on the manufacturer's priorities. For instance, a manufacturer may prefer to minimize false positives to detect and respond to attacks quickly or to investigate all suspicious cases. In our case, we set the value of ℓ to 200. However, the chosen threshold and ℓ value may have the unintended consequence of missing some low-intensity and short-duration attacks.

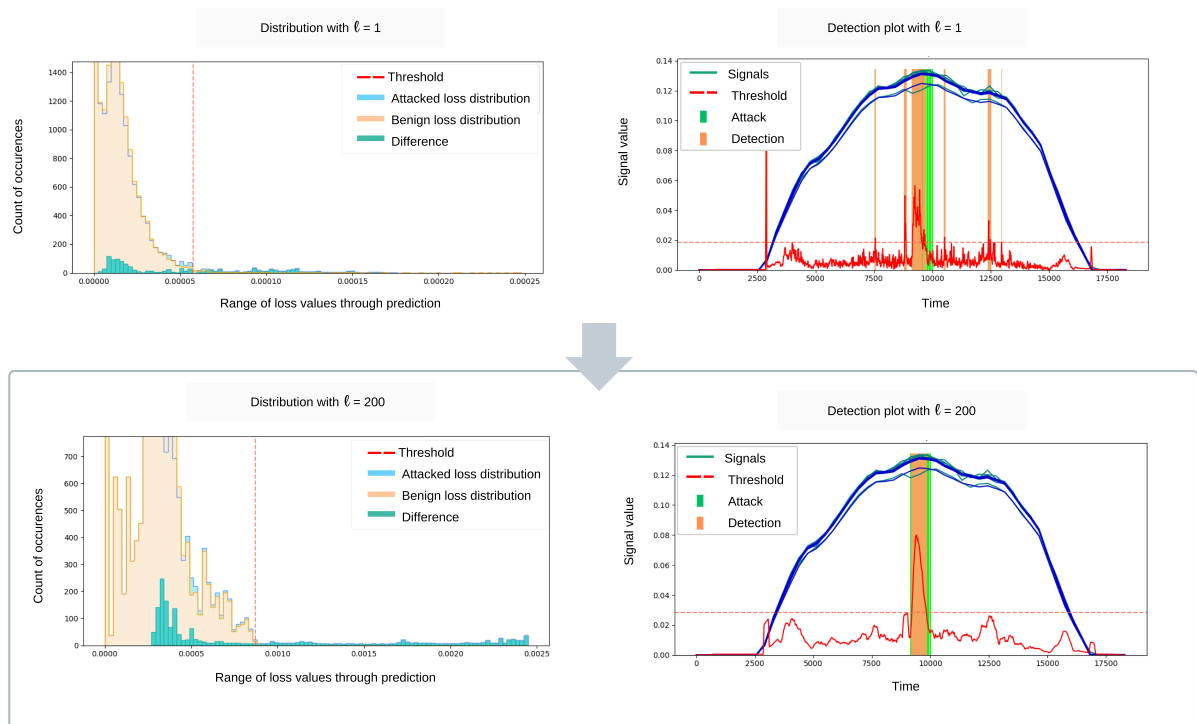


Figure 5.3: Visualization of the effect of different ℓ values, showing the distribution of the resulting loss values and a detection plot for each value ℓ .

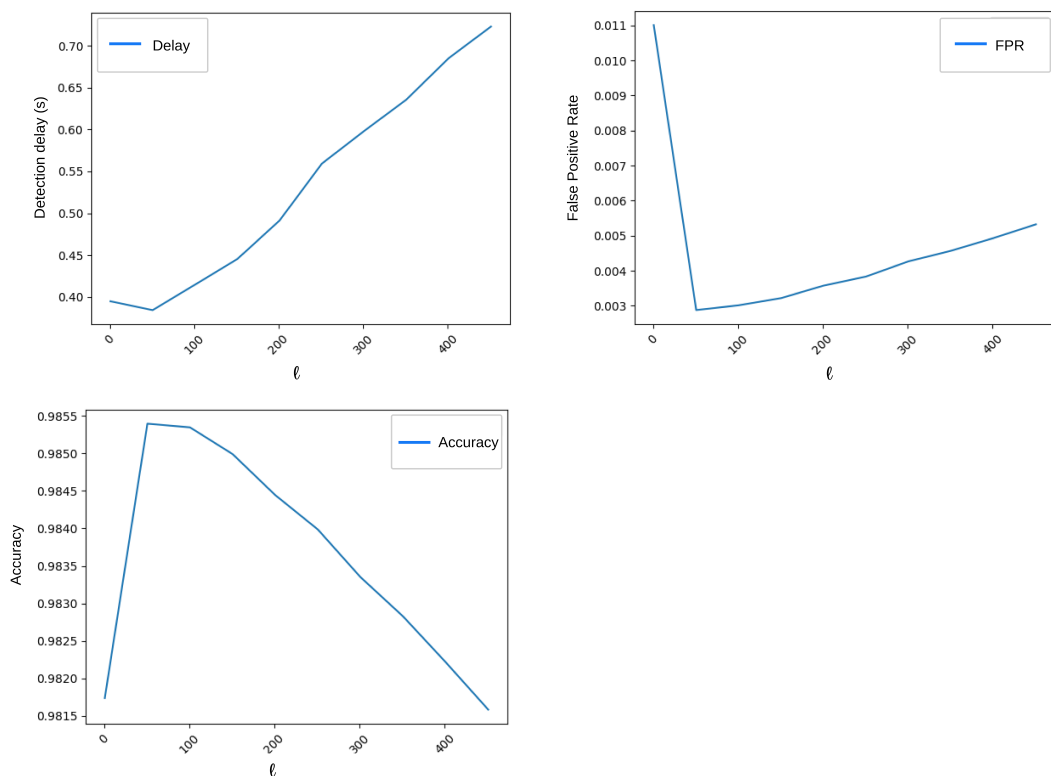


Figure 5.4: Visualization of the effect of different ℓ values, showing metrics (accuracy, detection delay, and false positive rate).

5.5 Discussion

5.5.1 Why Grouping Correlated Signals?

The joint forecasting of correlated signals offers several advantages for anomaly detection. First, it allows a single model per group to leverage the inherent interdependencies among group members, resulting in more accurate forecasts for each signal within the group. Second, any malicious modification of a signal is likely to impact the predictions of all group members, thereby increasing the cumulative prediction error as described in Eq. (5.2). This enhances the detectability of attacks compared to prior methods in the literature, as demonstrated in Chapter 6. Finally, instead of creating a stand-alone model for each individual signal as in [6], our approach requires the construction of only K forecasting models, rendering it a more appealing choice in resource-constrained environments.

5.5.2 Cost Analysis, Scalability and Robustness

When designing a solution for the industry it is crucial to consider the limited resource availability of automotive embedded systems.

The cost of our approach is dominated by that of the forecasting models. Apart from the C forecasting models, $K \cdot w$ signal values are stored for forecasting and $K \cdot \ell$ error values for decision purposes. The forecasting models are trained off-line in parallel, and the trained models are deployed in the vehicle. Therefore, the computational cost is dominated by the inference time of the forecasting models, where the inference processes of models are parallelizable.

It's important to highlight that during the offline training phase before deployment, our TCN architecture doesn't demand an excessively large model size, making the training phase feasible even for an average manufacturer. Clustering of signals also plays a significant role in reducing the computational requirements, resulting in fewer models.

In summary, these arguments substantiate that our system is well-suited for deployment in industrial environments, especially in the context of automotive embedded systems where resource constraints are a critical consideration.

Chapter 6

Evaluation

6.1 Dataset

We use two CAN datasets for evaluation: Dataset-1 introduced in [6], and Dataset-2 introduced in [12]. Dataset-1 contains seven short (<1 minute) traces of specific driving and traffic scenarios, and a longer trace (~25 minutes). Dataset-2 contains nine short traces and eleven longer traces.

As the datasets originate from the same vehicle type, both have 20 message IDs and 1-6 signals per ID. Similarly, both datasets contain message injection and message modification attacks. As our objective is to detect modification attacks, we only use the corresponding traces.

We evaluate our mechanism on Dataset-1 to compare its performance to the chosen baseline described in Section 6.3. Since the two datasets are based on very similar CAN traffic from the same vehicle type, and most attacks follow the same strategy (only the RANDOM and DELTA attacks are not included in both), we present only the joint results.

The attacks have been performed using 6 different signal modification strategies:

- ADD-DECR - Modify with decrement value: a decrease per message is subtracted from the original value.
- ADD-INCR - Modify with increment: increases the original value by one increment per message.
- CONST - Change to constant: constant value replaces the original value.
- NEG-OFFSET - Modify with delta: a given value is subtracted from the original data value.
- POS-OFFSET - Modify with delta: a given value is added to the original data value.
- REPLAY - Replace the original data value with a previous value.
- DELTA - An attacker chosen value is added to the original value.
- RANDOM - The original value is replaced by a new random value in every attacked message.

For illustration, an example for a REPLAY modification attack is depicted in Figure 6.1.

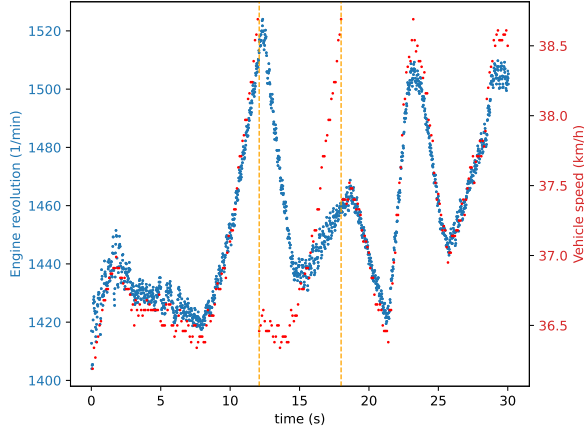


Figure 6.1: REPLAY attack, shown between the vertical lines, targeting messages with ID 0410, modifying speed signals [12].

6.2 Model Architecture and Parameters

For evaluation, we instantiate our proposal described in Chapter 5. We create two datasets for training and testing purposes. A total number of 3.2 million CAN messages were used to create a training dataset for signal forecasting and calibrating all parameters of our approach (i.e., K , C , w , ℓ). Our calibrated model is tested on 1.3 million benign and malicious test messages (67 attacked traces and 9 benign traces), each containing one attacked signal. Both datasets undergo the same pre-processing steps with the same parameters that were computed exclusively on the training data.

Pre-processing: We use a signal mask based on the bit flip rate to extract relevant signals. We retain $K = 20$ of the $N = 77$ extracted signals that describe the state of the vehicle and likely to have sufficient predictive power for signal forecasting¹. The retained signals are normalized as described in Section 5.1.

Signal grouping: We conduct a correlation analysis on the signals and identify groups of correlated signals. We utilize hierarchical clustering with Pearson correlation as a similarity measure, and group linearly dependent signals together accordingly. We identify $C = 9$ clusters of the 20 signals in our dataset.

Signal forecasting: For forecasting, we use multi-channel Temporal Convolutional Networks (TCN). We apply an input sliding window of size $w = 1750$, equivalent to roughly 3 seconds, and each TCN has a receptive field with the same size w . Each channel of the multi-channel model corresponds to an individual signal in the group. The output of the TCN layers is then forwarded to a fully connected linear layer which generates the prediction of the upcoming signal values. Each multichannel TCN layer has four dilatation layers with a logarithmic offset of 2 (1,2,4,8). The kernel size is fixed at 16. We train each forecasting model with Adam optimizer and MSE loss using early stopping. This forecasting module is illustrated in Figure 6.2.

The total size of all forecasting models, capable of handling all message IDs together in groups, is approximately 15 MB and contains 4.157 million parameters.

¹Note that this information is already known to a car manufacturer

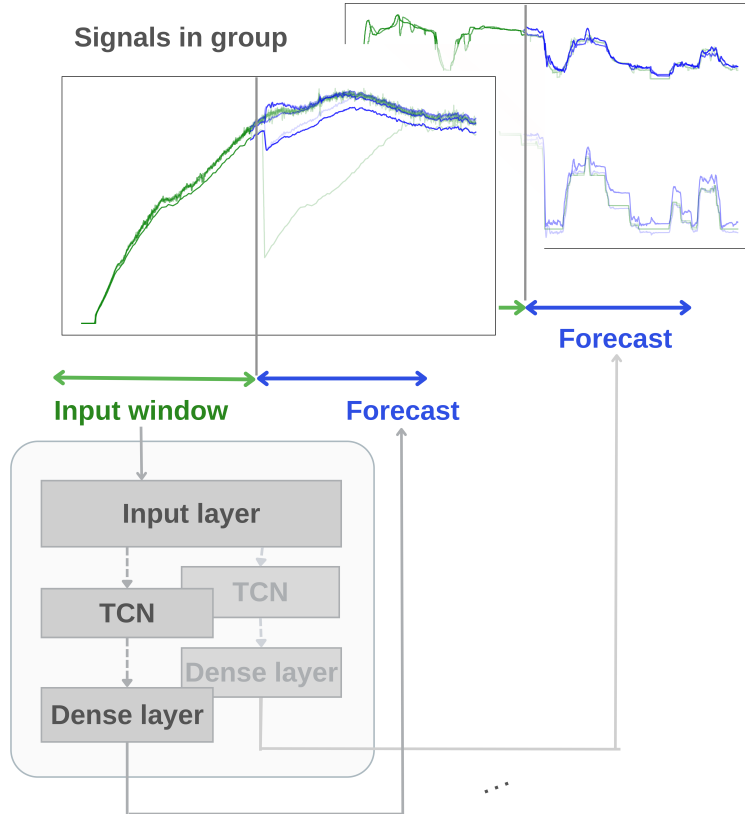


Figure 6.2: Visualization of the forecasting module.

Decision: We average the last $\ell = 200$ prediction error values of our forecasting models and compare with threshold τ which is calibrated according to the three-sigma rule on the training data as described in Section 5.4. In other words, we do *not* use the attacked traces in our dataset to adjust τ because it is unlikely to have sufficiently representative data about all possible attacks in practice. The decision module is illustrated in Figure 6.3.

6.3 Comparison with Baselines

The most relevant related works are CANet [16], INDRA [17], and the single TCN (S-TCN) anomaly detector architecture from [6]. To avoid confusion, from now on, we will refer to the Single TCN method (S-TCN), and refer to our proposed solution described in Section 6.2 as Correlation-based TCN (C-TCN).

The INDRA framework has been shown to outperform other relevant unsupervised approaches including CANet regarding false positives and detection accuracy. Moreover, according to numerical experiments on two datasets, the SynCAN dataset [16] and Dataset-1, the S-TCN approach has larger accuracy with a significantly lower false positive rate than INDRA. Therefore, it is sufficient to show that our solution outperforms the S-TCN approach, because it has demonstrated superior performance compared to CANet and INDRA [6].

To properly compare the two results, we adapt the S-TCN approach by training one TCN model per

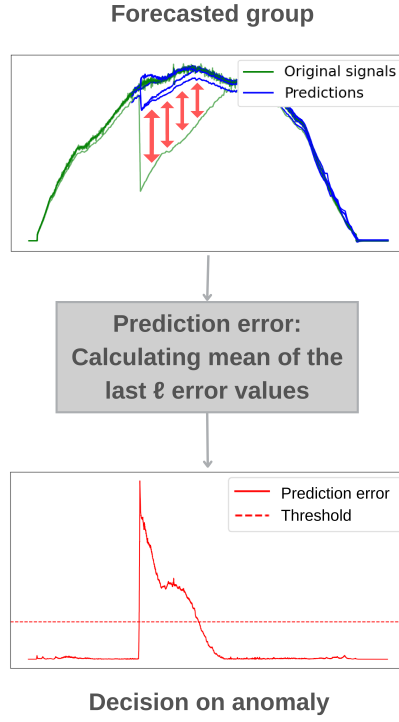


Figure 6.3: Visualization of the decision module.

signal but keeping the rest of the process, i.e., the data pre-processing, the same as our C-TCN solution. As expected, this adapted approach can reconstruct the expected behavior of CAN signals individually.

6.4 Evaluation Metrics

We evaluate both the baseline S-TCN and our proposed C-TCN method using standard performance metrics: accuracy, false positive rate, precision, and recall.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (6.1)$$

$$\text{Precision} = \frac{TP}{TP + FP} \quad (6.2)$$

$$\text{Recall} = \frac{TP}{TP + FN} \quad (6.3)$$

$$\text{FPR} = \frac{FP}{FP + TN} \quad (6.4)$$

Metrics are calculated according to Equations 6.1, 6.2, 6.3, 6.4, where TP means the number of true positive detections, FP the number of false positives, TN the true negatives and FN the false negatives.

Precision and recall are particularly important metrics in this context, since the *testing* dataset is often imbalanced; attacks on the CAN bus are often short, which means that the number of benign instances significantly exceeds the number of attack instances.

In addition, we also measure the time it takes to detect attacks (denoted by T_D), and the fraction of

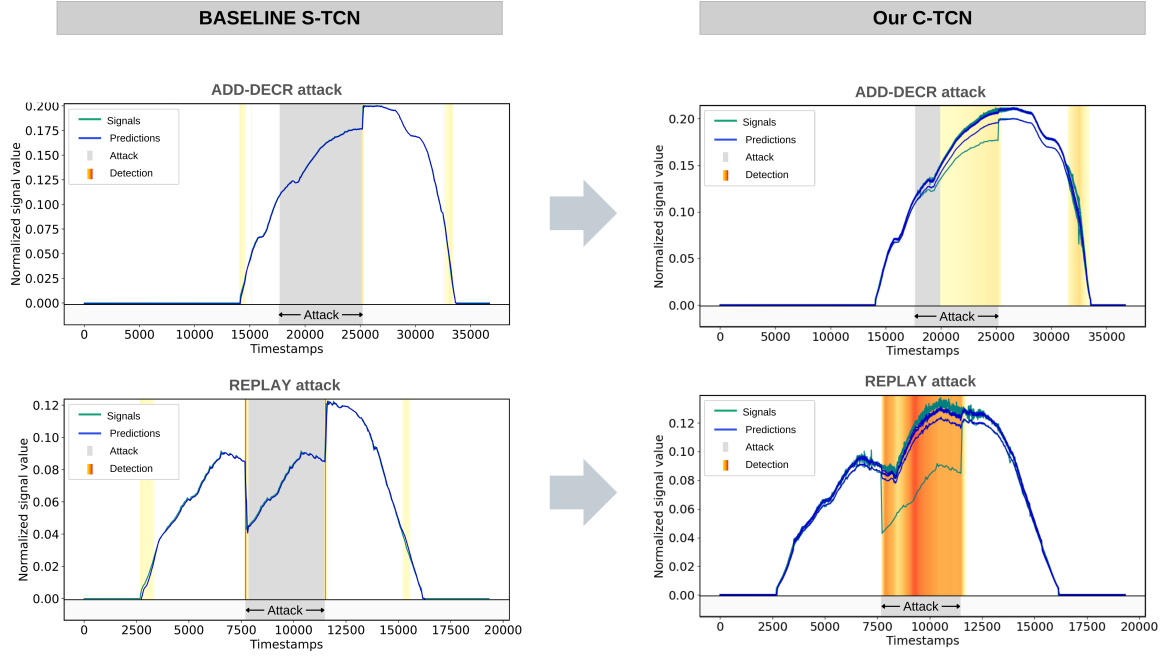


Figure 6.4: Comparative evaluation of S-TCN vs. C-TCN on two attacked traces. ADD-DECR (add decrement value) attack (first row of each column) and a REPLAY attack (second row of each column) are shown. The figure shows the attacked region marked by grey vertical lines and detections marked by yellow to red vertical lines, with the magnitude of the cumulative prediction error indicated by the darkness of the color.

attacked traces that are successfully detected (denoted by R_D):

$$T_D = \frac{\sum_{n=1}^{N_t} (t_{detection} - t_{attack})}{N_t} \quad (6.5)$$

$$R_D = \frac{\sum_{n=1}^{N_t} \mathbb{1}_{\{\text{trace } n \text{ is detected as anomalous}\}}}{N_t} \quad (6.6)$$

where N_t is the number of attacked traces, $t_{detection}$ is the time of detection (time of the first message whose signal values trigger anomaly), t_{attack} is the starting time of the attack (time of first attacked message) and $\mathbb{1}$ is the indicator function. Note that, while recall in Eq. (6.1) measures the detection performance on individual messages, detection rate measures the recall with respect to the traces. Indeed, both datasets used for evaluation includes short driving scenarios affected by various types of attacks, as described in Section 6.1, and an attacked trace is successfully detected if at least one message belonging to the attacked section of the trace triggers detection.

6.5 Results

All experiments were done using the TCN implementation in Keras [24].

Table 6.3 shows the accuracy and false positive rate for benign and malicious test sets, as well as the precision, recall, detection rate, and detection delay for attacked traces for both message modification and

Table 6.1: Comparing detailed results of evaluating the baseline S-TCN and the proposed correlation-based C-TCN on each modification attack types from both dataset.

	Model	Accuracy	FPR	Precision	Recall	R_D
ADD-DECR	S-TCN	0.93	0.06	0.35	0.34	0.45
	C-TCN	0.97	0.05	0.78	0.73	0.95
ADD-INCR	S-TCN	0.91	0.05	0.28	0.13	0.43
	C-TCN	0.97	0.04	0.79	0.71	0.96
CONST	S-TCN	0.91	0.04	0.09	0.01	0.0
	C-TCN	0.97	0.04	0.64	0.62	0.8
NEG-OFFSET	S-TCN	0.94	0.02	0.18	0.04	1.00
	C-TCN	0.98	0.02	0.75	1.00	1.00
POS-OFFSET	S-TCN	0.94	0.02	0.18	0.04	1.00
	C-TCN	0.98	0.02	0.76	1.00	1.00
REPLAY	S-TCN	0.93	0.03	0.08	0.03	0.55
	C-TCN	0.96	0.02	0.80	0.65	1.00
DELTA	S-TCN	0.9	0.03	0.05	0.01	1.00
	C-TCN	0.99	0.06	0.86	0.87	0.88
RANDOM	S-TCN	0.97	0.11	0.84	0.99	1.00
	C-TCN	0.99	0.06	0.92	1.00	1.00

Table 6.2: Comparing detailed results of evaluating the baseline S-TCN and the proposed correlation-based C-TCN on each injection attack types.

	Model	Accuracy	FPR	Precision	Recall	R_D
ADD-DECR INJ	S-TCN	0.95	0.01	0.54	0.11	0.67
	C-TCN	0.98	0.01	0.91	0.59	1.00
ADD-INCR INJ	S-TCN	0.95	0.01	0.59	0.11	0.70
	C-TCN	0.98	0.00	0.77	0.42	0.80
CONST INJ	S-TCN	0.95	0.01	0.51	0.09	0.60
	C-TCN	0.98	0.00	0.79	0.42	0.83
NEG-OFFSET INJ	S-TCN	0.97	0.02	0.85	0.61	1.00
	C-TCN	1.00	0.01	0.92	1.00	1.00
POS-OFFSET INJ	S-TCN	0.97	0.02	0.84	0.57	1.00
	C-TCN	1.00	0.01	0.92	0.99	1.00
REPLAY INJ	S-TCN	0.96	0.01	0.69	0.17	1.00
	C-TCN	0.99	0.01	0.97	0.79	1.00

message injection attacks. These metrics are calculated across multiple traces and averaged to provide the overall results shown in the table.

All metrics are also calculated for each attack type individually to determine the effectiveness against each type. Detailed results for message modification attacks in Table 6.1 show that both our solution and

Table 6.3: Comparing overall results of evaluating the baseline S-TCN and the proposed correlation-based C-TCN on benign and malicious test traces from both dataset.

	BENIGN		MALICIOUS MODIFICATION		MALICIOUS INJECTION	
	S-TCN	C-TCN	S-TCN	C-TCN	S-TCN	C-TCN
Accuracy	0.98	0.99	0.93	0.98	0.96	0.99
FPR	0.03	0.02	0.05	0.04	0.01	0.01
Precision	-	-	0.30	0.80	0.67	0.88
Recall	-	-	0.24	0.83	0.28	0.70
R_D	-	-	0.68	0.95	0.79	0.94

the baseline solution can easily detect modification attacks like NEG-OFFSET, POS-OFFSET, RANDOM, DELTA and some REPLAY attack. The baseline S-TCN performs poorly against the stealthier ADD-DECR, ADD-INCR and CONST attacks, while our results are 95%, 96%, and 80% respectively.

To investigate the use of only one IDS system in a vehicle, we also tested our solution against message injection attacks. Although we do not focus on detecting these attacks, we demonstrate that the solution can be applied to detect message injections as well. The results for message injection attacks in Table 6.2 show that, on average, both S-TCN and C-TCN perform better in this case than in the case of message modification attacks because duplicate signal values result in a more significant anomaly, as the values frequently oscillate between the two values.

After experimenting, we conclude that correlation-based C-TCN can effectively detect attacks on CAN bus data. Our major findings are as follows:

1. Grouping of CAN signals based on correlation improves the detection performance from 68% to 95% which means that our proposed C-TCN method can detect 95% of all the modification attack scenarios. These attacks are detected with a delay of 0.38 seconds on average.
2. Correlation-based C-TCN significantly outperforms S-TCN on all evaluated metrics, especially regarding precision and recall, where C-TCN achieves 80-83% average performance.
3. Table 6.1 shows that our C-TCN can detect even stealthier attacks that do not significantly modify signals (i.e. ADD-DECR, ADD-INCR and CONST attacks). Figure 6.4 presents an example of this improvement over the S-TCN baseline.
4. In addition to modification attacks, C-TCN also effectively identifies injection attacks, allowing the identification of both types of attacks by a single algorithm.

As Figure 6.4 shows, S-TCN fails to detect the stealthier ADD-DECR attack, which slowly modifies the original signal message-by-message. It is only detected when the attack abruptly stops, and the signal returns to its original value. In contrast, our C-TCN model can detect the attack earlier when the modification induces a detectable change in the cumulative prediction error. Similarly, while both models can detect the start of a replay attack, the baseline S-TCN cannot detect it throughout the entire attack span, whereas our C-TCN can.

Chapter 7

Conclusion

This paper presented a novel approach to intrusion detection on the CAN bus. We aimed at detecting message modification attacks, the most complex attack type possible on the CAN bus. We showed that a correlation-based TCN model can efficiently predict the subsequent values of the vehicle signals, which can be used for anomaly detection. Finally, we also presented measurements demonstrating that our approach outperforms the state-of-the-art.

Our main contribution is to combine correlation analysis with time-series forecasting to improve detection accuracy. By grouping signals first based on their correlation, we create models that can predict future values with a high accuracy. During an attack, the forecasting of a group of correlated signals is significantly less accurate, allowing the detection of the anomaly. Furthermore, by grouping the signals, we can use fewer models resulting in a smaller footprint, which is an important factor for embedded systems.

In case an attacker knows which signals are clustered together and understands how the signals usually behave, it may be able to modify all the signals in the group without being detected. This requires maintaining the normal signal behavior including the inter-dependencies between different signals. However, it is unlikely that the attacker have all these capabilities in practice, especially if the groups are sufficiently large and the device running our integrated solution is adequately protected.

In our future work, we plan to analyze correlations in different traffic situations to improve our solution.

Acknowledgements

I would like to express my gratitude to my mentors, András Gazdag, Gergely Ács, and Gergely Biczók, for their invaluable assistance and guidance.

Bibliography

- [1] Omid Avatefipour and Hafiz Malik. State-of-the-art survey on in-vehicle network communication (can-bus) security and vulnerabilities, 2018.
- [2] Stefan Axelsson. Intrusion detection systems: A survey and taxonomy. Technical report, Stockholm University, 04 2000.
- [3] Mehmet Bozdal, Mohammad Samie, Sohaib Aslam, and Ian Jennions. Evaluation of can bus security challenges. *Sensors*, 20(8), 2020. ISSN 1424-8220. doi: 10.3390/s20082364. URL <https://www.mdpi.com/1424-8220/20/8/2364>.
- [4] Miguel Á. Carreira-Perpiñán. A review of mean-shift algorithms for clustering, 2015.
- [5] Stephen Checkoway, Damon McCoy, Brian Kantor, Danny Anderson, Hovav Shacham, Stefan Savage, Karl Koscher, Alexei Czeskis, Franziska Roesner, and Tadayoshi Kohno. Comprehensive experimental analyses of automotive attack surfaces. In *20th USENIX Security Symposium (USENIX Security 11)*, San Francisco, CA, aug 2011. USENIX Association. URL <https://www.usenix.org/conference/usenix-security-11/comprehensive-experimental-analyses-automotive-attack-surfaces>.
- [6] Irina Chiscop, András Gazdag, Joost Bosman, and Gergely Biczók. Detecting message modification attacks on the CAN bus with temporal convolutional networks. In *Proceedings of the 7th International Conference on Vehicle Technology and Intelligent Transport Systems*, 2021. doi: 10.5220/0010445504880496. URL <https://doi.org/10.5220/0010445504880496>.
- [7] Kyong-Tak Cho and Kang G. Shin. Fingerprinting electronic control units for vehicle intrusion detection. In *Proceedings of the 25th USENIX Security Symposium*, pages 911–927. USENIX Association, 2016. ISBN 9781931971324. URL https://www.usenix.org/system/files/conference/usenixsecurity16/sec16_paper_cho.pdf.
- [8] Mohamed-Cherif Dani, François-Xavier Jollois, Mohamed Nadif, and Cassiano Freixo. Adaptive threshold for anomaly detection using time series segmentation. In Sabri Arik, Tingwen Huang, Weng Kin Lai, and Qingshan Liu, editors, *Neural Information Processing*, pages 82–89, Cham, 2015. Springer International Publishing. ISBN 978-3-319-26555-1.
- [9] Essam F. El-Hashash and Raga Hassan Ali Shiekh. A comparison of the pearson, spearman rank and kendall tau correlation coefficients using quantitative variables. *Asian Journal of Probability and Statistics*, 20(3):36–48, Oct. 2022. doi: 10.9734/ajpas/2022/v20i3425. URL <https://journalajpas.com/index.php/AJPAS/article/view/425>.
- [10] Absalom E. Ezugwu, Abiodun M. Ikotun, Olaide O. Oyelade, Laith Abualigah, Jeffery O. Agushaka, Christopher I. Eke, and Andronicus A. Akinyelu. A comprehensive survey of clustering algorithms: State-of-the-art machine learning applications, taxonomy, challenges, and future research prospects. *Engineering Applications of Artificial Intelligence*, 110:104743, 2022.

ISSN 0952-1976. doi: <https://doi.org/10.1016/j.engappai.2022.104743>. URL <https://www.sciencedirect.com/science/article/pii/S095219762200046X>.

- [11] Brendan J. Frey and Delbert Dueck. Clustering by passing messages between data points. *Science*, 315(5814):972–976, 2007. doi: 10.1126/science.1136800. URL <https://www.science.org/doi/abs/10.1126/science.1136800>.
- [12] András Gazdag, Rudolf Ferenc, and Levente Buttyán. Crysyst dataset of can traffic logs containing fabrication and masquerade attacks. *Scientific Data*, 2023.
- [13] András Gazdag, Dóra Neubrandt, Levente Buttyán, and Zsolt Szalay. Detection of injection attacks in compressed can traffic logs. In *International Workshop on Cyber Security for Intelligent Transportation Systems, Held in Conjunction with ESORICS 2018*. Springer, 2018.
- [14] András Gazdag, György Lupták, and Levente Buttyán. Correlation-based anomaly detection for the can bus. In *Euro-CYBERSEC, Nice, France, 2021*.
- [15] Jake Guidry, Fahad Sohrab, Raju Gottumukkala, Satya Katragadda, and Moncef Gabbouj. One-class classification for intrusion detection on vehicular networks, 2023.
- [16] Markus Hanselmann, Thilo Strauss, Katharina Dormann, and Holger Ulmer. Canet: An unsupervised intrusion detection system for high dimensional can bus data. *IEEE Access*, 8:58194–58205, 2020. doi: 10.1109/ACCESS.2020.2982544.
- [17] Vipin Kumar Kukkala, Sooryaa Vignesh Thiruloga, and Sudeep Pasricha. Indra: Intrusion detection using recurrent autoencoders in automotive embedded systems, 2020.
- [18] Seyoung Lee, Hyo Jin Jo, Aram Cho, Dong Hoon Lee, and Wonsuk Choi. Ttids: Transmission-resuming time-based intrusion detection system for controller area network (can). *IEEE Access*, 10:52139–52153, 2022. doi: 10.1109/ACCESS.2022.3174356.
- [19] Mirco Marchetti and Dario Stabili. Read: Reverse engineering of automotive data frames. *IEEE Transactions on Information Forensics and Security*, 14(4):1083–1097, 2019. doi: 10.1109/TIFS.2018.2870826.
- [20] Moti Markovitz and Avishai Wool. Field classification, modeling and anomaly detection in unknown can bus networks. *Vehicular Communications*, 9, 03 2017. doi: 10.1016/j.vehcom.2017.02.005.
- [21] Pablo Moriano, Robert A. Bridges, and Michael D. Iannacone. Detecting can masquerade attacks with signal clustering similarity. *ArXiv*, abs/2201.02665, 2022. URL <https://api.semanticscholar.org/CorpusID:245836760>.
- [22] Michael Müter and Naim Asaj. Entropy-based anomaly detection for in-vehicle networks. *2011 IEEE Intelligent Vehicles Symposium (IV)*, pages 1110–1115, 2011. URL <https://api.semanticscholar.org/CorpusID:9017380>.
- [23] Brent C. Nolan, Scott Graham, Barry Mullins, and Christine Schubert Kabban. Unsupervised time series extraction from controller area network payloads. In *Proceedings of the 2018 IEEE 88th Vehicular Technology Conference (VTC-Fall)*, pages 1–5. IEEE, 8 2018. ISBN 978-1-5386-6358-5. doi: 10.1109/VTCFall.2018.8690615. URL <https://ieeexplore.ieee.org/document/8690615/>.

- [24] Philippe Remy. Temporal convolutional networks for keras. <https://github.com/philipperemy/keras-tcn>, 2020.
- [25] Mayra Z. Rodriguez, Cesar H. Comin, Dalcimar Casanova, Odemir M. Bruno, Diego R. Amancio, Luciano da F. Costa, and Francisco A. Rodrigues. Clustering algorithms: A comparative approach. *PLOS ONE*, 14:1–34, 01 2019. doi: 10.1371/journal.pone.0210236. URL <https://doi.org/10.1371/journal.pone.0210236>.
- [26] Hyun Min Song, Ha Rang Kim, and Huy Kang Kim. Intrusion detection system based on the analysis of time intervals of can messages for in-vehicle network. *2016 International Conference on Information Networking (ICOIN)*, pages 63–68, 2016. URL <https://api.semanticscholar.org/CorpusID:9333718>.
- [27] Miki E. Verma, Robert A. Bridges, Jordan J. Sosnowski, Samuel C. Hollifield, and Michael D. Iannacone. Can-d: A modular four-step pipeline for comprehensively decoding controller area network data. *IEEE Transactions on Vehicular Technology*, 70:9685–9700, 10 2021. ISSN 19399359. doi: 10.1109/TVT.2021.3092354.
- [28] Clinton Young, Habeeb Olufowobi, Gedare Bloom, and Joseph Zambreno. Automotive intrusion detection based on constant can message frequencies across vehicle driving modes. In *AutoSec '19: Proceedings of the ACM Workshop on Automotive Cybersecurity*, pages 9–14, 03 2019. doi: 10.1145/3309171.3309179.

Appendix



Figure A.1: Clustering results for Affinity Propagation, DBSCAN, Hierarchical and MeanShift clustering (from top to bottom) calculated with the Kendall, Spearman and Pearson correlation (from left to right).