

# Detecting Malicious Data Injections in Wireless Sensor Networks: A Survey

VITTORIO P. ILLIANO and EMIL C. LUPU, Imperial College London

Wireless Sensor Networks are widely advocated to monitor environmental parameters, structural integrity of the built environment and use of urban spaces, services and utilities. However, embedded sensors are vulnerable to compromise by external actors through malware but also through their wireless and physical interfaces. Compromised sensors can be made to report false measurements with the aim to produce inappropriate and potentially dangerous responses. Such malicious data injections can be particularly difficult to detect if multiple sensors have been compromised as they could emulate plausible sensor behaviour such as failures or detection of events where none occur. This survey reviews the related work on malicious data injection in wireless sensor networks, derives general principles and a classification of approaches within this domain, compares related studies and identifies areas that require further investigation.

Categories and Subject Descriptors: K.6.5 [Management of Computing and Information Systems]: Security and Protection—*Unauthorised access*

General Terms: Security, Algorithms, Measurement

Additional Key Words and Phrases: Wireless sensor networks, security, correlation

## ACM Reference Format:

Vittorio P. Illiano and Emil C. Lupu. 2015. Detecting malicious data injections in wireless sensor networks: A survey. *ACM Comput. Surv.* 48, 2, Article 24 (October 2015), 33 pages.

DOI: <http://dx.doi.org/10.1145/2818184>

## 1. INTRODUCTION

Wireless sensor networks (WSNs) are an attractive solution to the problem of collecting data from physical spaces, thanks to their flexibility, low cost, and ease of deployment. Applications of WSNs include a broad variety of tasks in both shared and personal environments. In shared environments, applications include monitoring of infrastructures such as the water network, improvement of road traffic, monitoring of environmental parameters and surveillance. In personal environments, applications include monitoring homes for energy efficiency, user activity such as exercise and sleep, and physiological parameters for health care through both wearable and implantable sensors.

In some aspects, WSNs are similar to traditional wired and wireless networks, but they also differ in other aspects, such as the sensors' limited computational and power resources. Sensors need to be cheap, be physically small, communicate wirelessly, and have low-power consumption whether to monitor a human body or a large flood plain, and therein lie their main advantages. But these characteristics are also their main

---

Authors' addresses: V. P. Illiano and E. C. Lupu, Computer Science Department, Imperial College London, Department of Computing, Huxley Building 180 Queen's Gate, South Kensington Campus, London SW7 2AZ, UK; emails: {v.illiano13, e.c.lupu}@imperial.ac.uk.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or [permissions@acm.org](mailto:permissions@acm.org).

© 2015 ACM 0360-0300/2015/10-ART24 \$15.00

DOI: <http://dx.doi.org/10.1145/2818184>

limitations as they lead to more frequent failures, poor physical protection, limited degree of redundancy and processing, and limited ability to carry out complex operations.

Wireless sensors carry a much higher risk of being compromised. Their deployments are often unattended and physically accessible, and use of tamper-resistant hardware is often too expensive. The wireless medium is difficult to secure and can be compromised at all layers of the protocol stack. Cryptographic operations and key management consume valuable computational and power resources and cannot provide a solution once a node has been compromised. Yet, despite this, WSNs are increasingly used to monitor critical infrastructures and human health where malicious attacks can lead to significant damage and even loss of life.

Faced with the challenge of securing WSNs, researchers have proposed new security solutions for these platforms. The literature is rich and we can only cite a few examples [Karlof and Wagner 2003; Perrig et al. 2004; Du et al. 2005; Liu and Ning 2008; Khan and Alghathbar 2010]. Most studies focus on proposing solutions against physical-level and network-level threats, such as jamming attacks, attacks against the routing protocols, confidentiality, and integrity of the data in transit. Another body of work is that of *software attestation*, which assesses the node integrity and in particular checks that the nodes run the expected software [Seshadri et al. 2004; Park and Shin 2005; Seshadri et al. 2006; Zhang and Liu 2010].

Despite such solutions, many attacks remain possible against wireless sensor nodes. For example, an attacker may compromise a node through its physical interfaces or tamper with the node hardware itself in order to introduce wrong measurements in the network. This would defeat many of the solutions presented in the literature as the cryptographic material present on a compromised sensor would (in the absence of trusted hardware) be available to the attacker. Even when the sensors are hard to reach or to tamper with, an attacker may also seek to compromise the measurements by locally manipulating the sensed environment to induce malicious readings, for example, using a lighter to trigger a fire alarm. We refer to all these kinds of attacks as *malicious data injections*. Their aim is to compromise the mission of the WSN by producing a picture about the sensed phenomenon, which is different from the real one with potentially devastating effects. In particular, an attacker may seek to

- elicit an inappropriate system response**, for example, triggering an overload on a power grid, leading to partial shutdown; or
- masking a desired system response**, for example, silencing an intrusion alarm.

Protecting from such attacks becomes essential because of their potential impact, and this survey focuses on solutions proposed that could address this problem. The main challenge for detecting malicious data injections is finding sufficient evidence of the attack. A possible approach is to look for evidence of tampering with the sensor itself through software attestation, as mentioned earlier. However, software attestation is difficult to deploy in practice (e.g., because of timeliness constraints and device hardware restrictions [Castelluccia et al. 2009]). Attacks that locally modify the sensed environment are also still possible. Another approach is to look for evidence of changed traffic patterns in the communication between the sensors, for example, through traffic analysis [Buttyan and Hubaux 2008]. While effective for detecting network-level attacks, in particular on routing, such approaches often cannot detect malicious data injections since an attacker may modify the values reported by the sensors without changing the traffic patterns of the communications between sensors.

For these reasons, we focus in this article on techniques that look for evidence of compromise in the sensor measurements themselves, *regardless of how they may have been compromised*. Thus, we include in the scope of this survey techniques that perform data analyses on such measurements to detect malicious interference. In addition, we

include papers that aim to detect generic anomalies in WSNs but that are still based on the collected measurements. In contrast, anomaly-based techniques that operate on network parameters such as packet transmission rate, packet drop rate, transmission power, and so forth are beyond the scope of this survey. Indeed, a key aspect of the detection of malicious data injections is the construction of the *data expectation model*, that is, the model that allows one to define expectations about the sensors' measurements. In this context, anomalies arise in the correlation structures that are natively present in the data itself, which cannot be found in network parameters, and may occur without any disruption to the network parameters.

All the papers reviewed in this survey assume that the attacker aims to cause noticeable undesired effects and injects measurements that differ in some detectable way from the correct values that should be reported at that point in time and space. This is the assumption that enables the use of data analysis to detect data injections. However, note that the real value that *should* be reported by compromised sensors is not observable directly. Instead, it can only be characterized from indirect information such as values reported by other sensors, which may or may not be sufficient to detect the compromise. The problem is even more difficult as the indirect information may itself not be correct due to the presence of faults or naturally occurring events. *Faults* refer to any kind of genuine errors, transient or not, and may be difficult to distinguish from a malicious injection. *Events* refer to substantial changes in the sensed phenomenon like a fire, an earthquake, and so forth. We refer to the problem of distinguishing malicious data injections from events and faults as *diagnosis* and review the state-of-the-art approaches to the problem. Another cause for unreliable indirect information is the presence of *colluding sensors*, that is, when multiple compromised sensors produce malicious values in a coordinated fashion. In such scenarios, the attacker's leverage on the system increases and opens the possibility to new and more effective attacks.

Detecting and diagnosing malicious data injections is a subset of the more general problem of ensuring the integrity of the sensed data, which may have been corrupted by failures or in other ways. This is reflected in the studies surveyed, where many techniques designed, for example, for detecting faulty sensors or faulty data are also advocated for malicious data injections. Comparatively, only a small proportion of the papers explicitly focus on malicious data injections. However, there is a significant difference between faults and maliciously injected data since the latter is *deliberately* created in sophisticated ways to be difficult to detect. Therefore, there is a need for a survey that (1) analyzes the achievements and shortcomings of the work targeted to malicious data injections and (2) also reviews the state-of-the-art techniques proposed for nonmalicious data compromise and evaluates their suitability to this problem.

Within the context of WSNs, the applicable state-of-the-art studies broadly follow two types of approaches: *anomaly detection* techniques starting from about 2005 [Tanachaiwiwat and Helmy 2005] and *trust management* techniques starting from about 2006 [Zhang et al. 2006]. We review the state of the art for both approaches and compare the studies surveyed according to their

- adopted approach,
- ability to detect malicious data injections, and
- results and performance.

The remainder of this article is organized as follows. In Section 2, we describe existing surveys related to the one we present here. In Section 3, we recap concepts useful for understanding the rest of the article. In Section 4, we analyze possible ways of defining an expected behavior for sensor measurements and analyze the different approaches adopted in the state-of-the-art techniques. In Section 5, we analyze the state-of-the-art detection algorithms. In Section 6, we describe two aspects that are important to

tackle malicious data injections beyond detection: diagnosis and characterization of the attack. In Section 7, we give comparison tables for the techniques surveyed and their experimental results, together with a brief discussion. Finally, in Section 8, we present our conclusions and the open issues that emerged from this study.

## 2. RELATED SURVEYS

To the best of our knowledge, there are no previous surveys of techniques to detect malicious data injections in WSNs. Several surveys are, however, related, and we discuss them in this section.

Boukerche et al. [2008] analyze techniques for secure localization algorithms in WSNs. There are some similarities between malicious data injections and attacks on localization systems, since the sensors' location can be regarded as a particular physical phenomenon being sensed. However, many aspects of the techniques described in Boukerche et al. [2008] are specific to the localization problem. In particular, constraints on the topology, the radio transmission power, and delay provide a clear criterion to check the consistency of the information reported by the sensors. In contrast, we focus on techniques that do not require a priori knowledge of the physical phenomena monitored to check data consistency but examine and infer correlations from the data itself.

Rajasegarar et al. [2008] review 11 state-of-the-art papers about anomaly detection in WSNs. Although they focus on detecting intrusions, the survey also covers eliminating erroneous readings and reducing power consumption. The detection algorithms surveyed consider sensor measurements as well as network traffic and power consumption. In contrast, we focus on a more specific target: the detection of malicious data injections. We cover a broader spectrum of papers since we include techniques other than anomaly detection, describe further steps for detecting malicious data, and include a significant amount of literature published since then.

Xie et al. [2011] survey anomaly detection in WSNs, with a focus on the WSN architecture (Hierarchical/Flat) and the detection approach (statistical, rule based, data mining, etc.). They describe the detection procedure in a similar way to us: definition of a "normal profile," which we refer to as normal or expected behavior, and test to decide whether it is an anomaly or not, and to what extent. However, our survey is structured based on the approach to both the definition of the normal behavior and the detection based on it, while Xie et al. [2011] focus only on the latter. This choice allows us to pinpoint the motivation for the use of a particular detection technique, based on how the data normally looks. Moreover, the diagnosis process that classifies an anomaly as an attack is not analyzed in Xie et al. [2011], whereas it forms an important part of this survey.

Several surveys discuss trust management for security in WSNs (e.g., Lopez et al. [2010], Özdemir and Xiao [2009], and Sang et al. [2006]). However, they focus on attacks conducted through the network layer, while malicious data injections are given little attention. Yu et al. [2012] list all the threats that can be mitigated by trust management, including "stealthy attacks"—a kind of malicious data injection—but these are not analyzed in detail. Similarly, Zahariadis et al. [2010a] build a taxonomy of trust metrics, which includes consistency of reported values/data, but they focus mostly on the other network-related metrics. Also, Shen et al. [2011] survey defensive strategies against attacks to the network layer. In particular, such strategies are derived from game theory and take into account the strategies that can be adopted by the attacker to balance the profit and loss of reputation coming from the attack; in our survey, instead, we focus on techniques to assign and maintain such reputation.

The closest survey to the one presented here is Jurdak et al. [2011]. It describes anomaly detection strategies for detecting faults due to environmental factors (e.g.,

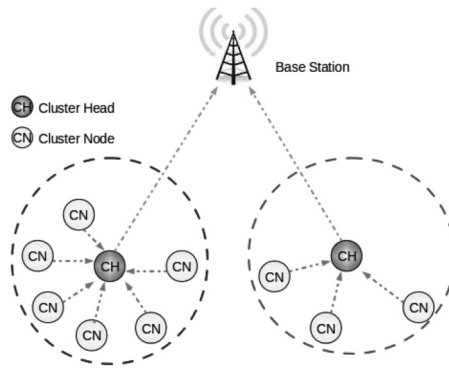


Fig. 1. LEACH measurements collection architecture.

obstructions near the sensor) or node hardware/software. Their description of anomaly detection is similar to ours, but the two surveys differ notably in the nature of the anomalies considered: attacks in our case, faults in theirs. Jurdak et al. [2011] also claim that anomalies can be detected by spatial or temporal comparisons between sensors, since it is unlikely that many sensors will exhibit a calibration skew or failure at the same time (assuming there are no group failures). This assumption considers anomalies (faults) as independent but does not hold in the presence of malicious data injections, in particular when there is collusion between the compromised sensors.

### 3. PRELIMINARIES

We describe in the following how sensors measurements are generally gathered in a WSN. We also introduce the two approaches used to detect malicious data injections so far: *anomaly detection* and *trust management*.

#### 3.1. Data Aggregation Schemes and Their Consequences

The typical workflow of a WSN starts with measuring a physical phenomenon through sensing devices connected to a wireless node that propagates the measurements through the network toward data sinks. Measurements collected and aggregated by data sinks (e.g., basestations) can then be interpreted or transmitted to a remote station. However, data can also be aggregated in the network by the intermediate transmitting nodes, with many possible variations on the aggregation architecture. The choice between the different schemes is based on criteria that optimize power efficiency, number of devices, coverage of the physical space, and so forth. Finding the optimal architecture based on such criteria remains an important research challenge.

Early work considered that all raw measurements are collected at the basestation, which performs data fusion and other computations [Shepard 1996; Singh et al. 1998]. Later on, especially after the introduction of the LEACH protocol [Heinzelman et al. 2000], architectures became increasingly hierarchical. LEACH applies a one-level hierarchy where sensors are organized in clusters and communicate with the cluster-head, which, in turn, communicates with the basestation, as shown in Figure 1. Cluster-based protocols, and especially those where the clusters change dynamically in time [Heinzelman et al. 2000], have proven to be more energy efficient when communication with the basestation requires multihop transmissions [Heinzelman et al. 2000].

The one-level hierarchy introduced in LEACH can be generalized to tree-based structures as described in Fasolo et al. [2007]. Intermediate tree nodes may simply merge the packets generated by different sources into a single packet without processing the data. This is referred to as *in-network aggregation without size reduction* [Fasolo et al.

2007]. Alternatively, they process the sensor measurements by applying aggregation operators (e.g., mean, minimum, maximum), which is referred to as *in-network aggregation with size reduction* [Fasolo et al. 2007]. So cluster heads assume the burden of the additional computation in order to minimize the data transmitted. In essence, this trades the power costs of computation for those of communication, but since in WSNs communication consumes much more power, the trade is usually favorable.

Information about the WSN architecture and where data aggregation is carried out is important for allocating the detection task to the WSN nodes. For instance, if in-network aggregation with size reduction is used, the basestation cannot analyze all the measurements and the aggregating nodes must assist the basestation in the detection task. In this case, the integrity of the aggregation process at these nodes must also be ascertained [Przydatek et al. 2003; Ganeriwal and Srivastava 2004; Roy et al. 2014].

### 3.2. Relationship to Anomaly Detection and Trust Management

Detection of malicious data injections has been addressed with two main approaches so far: *anomaly detection* (e.g., Tanachaiwiwat and Helmy [2005], Liu et al. [2007], and Sun et al. [2013]) and *trust management* (e.g., Atakli et al. [2008], Bao et al. [2012], and Oh et al. [2012]). While anomaly detection defines normal behaviors to infer the presence of anomalies, *trust management* evaluates the confidence level (trustworthiness) that a sensor's behavior is normal. Compromised sensors are then expected to get low trust values when they deviate from their expected behavior. Although anomaly detection is also based on the definition of an expected behavior—"anomaly detection refers to the problem of finding abnormalities in the data that do not conform to expected behaviour" Chandola et al. [2009]—the two approaches differ in how deviations are interpreted. In trust management, the sensors' measurements are analyzed with the granularity of a sensor, and each sensor has a trust value that is incrementally updated in time. Anomaly detection approaches, instead, can be applied with no restrictions in granularity from the single measurement to the whole system and generally work by defining a boundary for expected behavior such that everything outside that boundary is abnormal.

Given the similarities and differences between the two approaches, we structure the following two sections as follows: in the next section, we describe how to gather information about expected data, regardless of whether it is for anomaly detection or trust management. In Section 5, instead, we describe how to detect deviations from the expected data, treating anomaly detection and trust management separately.

## 4. MODELING EXPECTED DATA

In our context, *expected data* refers to a set of properties characterizing the measurements that are free of malicious injections. Given that no previous surveys focus on this issue, we start by introducing a generic formulation of WSN sensing. This enables us to analyze different models for the expected data and describe the related work with a coherent terminology as the terms used often differ from one article to another.

### 4.1. A Characterization of the Problem

We focus on interpreting the data and abstract from implementation-related issues, such as synchronization between sensors, and network-related issues, such as packet losses or delays. We consider a deployment region  $D$ , in which a set of  $N$  sensors are placed. Every sensor measures a physical attribute such as temperature, wind, water quality, power, and gas flows. The sensors' measurement process is characterized by a degree of uncertainty, which may be due to noise, faults, and malicious data injections. It is desirable to remove this uncertainty, so we introduce an ideal function  $\varphi$ , which represents the value of a sensor's reading in the absence of any source of uncertainty.

The independent variables of such function are the point in space  $s$  and the time instant  $t$  to which the readings correspond, as shown in Equation (1):

$$\varphi(s, t) \quad s \in D, \forall t. \quad (1)$$

We refer to this function as the *physical attribute function*. The reading produced at time  $t$  by a generic sensor  $i$ , deployed at position  $s_i$ , is some approximation of the physical attribute function evaluated at  $(s_i, t)$ . A generic sensor's reading can then be modeled as a function  $S_i$  that adds a generic measurement error  $\epsilon(s_i, t)$  to the physical attribute, which may change with time and space. Equation (2) defines the function  $S_i$  as follows:

$$S_i(t) = \varphi(s_i, t) + \epsilon(s_i, t) \quad i \in 1, \dots, N. \quad (2)$$

Note that the sensors' readings are the only observable quantities; both the physical attributes and the measurement errors are not observable directly. When malicious data injections occur, some of the sensors' readings also become unobservable, since the attacker substitutes fabricated measurements for the real ones. There is then the need to describe the real measurements with related information from some observable quantities. This process is effective if such related information allows us to discriminate injections and is itself not susceptible to injections.

Describing the unobservable real measurement in terms of observable properties is a modeling process that makes assumptions about how data can be described. For instance, the measurements produced by a sensor can be modeled as samples from a normal distribution [Zhang et al. 2006]. Assuming that compromised nodes do not produce data compliant with a normal distribution, the model can then discriminate compromised nodes [Zhang et al. 2006].

The relation that links the problem to a model is a one-to-many relation. Different models of the same problem are not equivalent, and choosing a good model is essential for good performance. In particular, a good model should be characterized by the following:

- Accuracy** – No model is perfect and every model is in fact an approximation. An accurate model minimizes the approximation error.
- Adaptability** – Physical attributes measured by the sensors change in time. As a consequence, models should adapt to the dynamically changing environment.
- Flexibility** – Good models should be applicable in a flexible way, regardless of the application. Such models should abstract as many details as possible and capture only those properties that are needed.

These desirable characteristics conflict with each other: accuracy may be better achieved with context-specific details, which limit flexibility and compromise adaptability. A particular adaptability requirement that significantly affects accuracy and flexibility is the sensors' *mobility*, as when sensor nodes migrate to new locations, previous expectations are invalidated. Indeed, sensor migrations correspond to a change in  $s_i$  in Equation (2), which potentially changes all the measurements' time series, leaving sensor-specific noise as the only invariant.

*Support for Mobility.* Even though mobility is an aspect that is not directly addressed in the detection of malicious data injections, some techniques are more suited to support mobile sensors than others. In particular, anomaly detection techniques that compare the measurements within a neighborhood without considering past behavior (e.g., Handschin et al. [1975], Ngai et al. [2006], Liu et al. [2007], Wu et al. [2007], and Guo et al. [2009]) can generally accommodate mobility, since for every time instant, new expectations are extracted. However, such techniques also need to become aware of topology changes in the presence of mobility.

Trust management techniques with exchanges of trust information (e.g., Bao et al. [2012], Huang et al. [2006], Ganeriwal et al. [2003], and Momani et al. [2008]) are also suited for mobility, since a sensor  $i$  that migrates to a new area and becomes a neighbor of  $j$  can benefit from recommendations from sensors that have been  $j$ 's neighbors in the past [Zahariadis et al. 2010b]. So far, exchanges of trust information have been considered without investigating the effects of mobility; therefore, sensor  $i$  will generally maintain indirect information about sensor  $j$  only if there is interaction between  $i$  and  $j$ , and  $i$  cannot observe  $j$ 's behavior (e.g., it is not in the wireless communication range). When sensors are mobile instead, even if  $i$  and  $j$  never interacted, they may interact in the future if they get closer. Only at that time, recommendations for  $j$  become of  $i$ 's interest, and a criterion to request such recommendations is needed.

The existing studies analyzed in the remainder of this work, by and large, ignore mobility aspects. We conclude, in light of the previous considerations, that more work is required to deal with the problems arising from the sensors' mobility.

## 4.2. Exploiting Correlation

Since the original measurements substituted with fabricated ones cannot be observed directly, they need to be characterized indirectly with related information. The relationship between two pieces of information is a *correlation*, which can be calculated online, with historical data, or modeled a priori. In either case, coexistence of genuine and compromised measurements may cause disruptions in the correlation, assuming that the correlations have not changed between the moment when they are calculated and the moment when they are used.

We refer here to *correlation* in a broad sense, meaning that there is some kind of continuous dependency, as opposed to Pearson's correlation coefficient, which is the most commonly used correlation metric. Referring to  $E$ ,  $\mu$ , and  $\sigma$  as the expected value, the mean, and the standard deviation, respectively, the Pearson correlation coefficient  $\rho_{XY}$  between two random variables  $X$  and  $Y$  is given by

$$\rho_{XY} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}. \quad (3)$$

Note that this coefficient measures only linear dependency between two variables, while nonlinear dependencies may be missed.

In WSNs, we can generally consider correlations across three different domains: *temporal*, *spatial*, and *attribute* domains [Rassam et al. 2013].

- Temporal correlation** is the dependency of a sensor's reading on its previous readings. It models the coherence in time of the sensed physical process.
- Spatial correlation** is the dependency in readings from different sensors at the same time. It models the similarities in how the sensed phenomenon is perceived by different sensors.
- Attribute correlation** is the dependency in readings that are related to different physical processes. It models physical dependencies among heterogeneous physical quantities such as temperature and relative humidity.

Usually a combination of these different kinds of correlation is used. We now analyze how they contribute to the definition of expected data.

## 4.3. Temporal Correlation

Variations in time of the sensed data can be modeled as a random process [Boukerche 2009], where the random variables at different times are correlated. As Equation (2) shows, the variation of a sensor's measurements in time depends on both the variations introduced by the physical attribute and the measurements' error. The variation of the



physical attribute in time is subject to constraints, such as the presence of gradual changes or the alternation of some patterns, since the phenomenon observed usually follows the laws of physics. So, if the measurements are gathered with sufficiently high frequency, consecutive measurements would be subject to similar constraints. This simple observation justifies a procedure that identifies errors (including malicious injections) when temporal variations do not respect these constraints. However, there are two main difficulties in applying this observation to assess deviations: the time evolution of the process is subject to uncertainty factors and the measurements are subject to noise.

When using Kalman filters [Kalman 1960] to model time series, these two factors are known respectively as *process noise* and *measurement noise*. The measurement noise is typically modeled as a Gaussian process. The process noise, instead, comes from the imperfections of the model used to describe the process dynamics. For example, when modeling the process as a discrete Markov process, the value at time  $t_1$  can be written as

$$\varphi(t_1) = F(\varphi(t_0)) + w(t_0), \quad (4)$$

where  $F$  models the expected evolution of the time process and  $w$  is the process noise.

The use of a Markovian process, modeled with a Kalman filter, forms the basis of the Extended Kalman Filter (EKF)-based algorithm by Sun et al. [2013]. Here, each sensor builds up a prediction for its neighbors as a function of the neighbors' previous reading. The difference between the predicted and the actual value forms a deviation that is used to detect malicious data injections. However, the authors point out that an attacker can elude the EKF algorithm by introducing changes that are sufficiently small. To address this shortcoming, the authors apply the CUSUM GLR algorithm, which considers the cumulative deviation across more time samples and tests it to be zero-mean. This property makes it more difficult for attackers to introduce deviations that achieve their goal.

Subramaniam et al. [2006] also define expected data with temporal correlation. Here, the authors fit the Probability Density Function (PDF) of the measurements inside a temporal window, through kernel density estimators. Given a new measurement  $p$ , the PDF gives information about the expected number of values falling in  $[p - r, p + r]$  (with parameter  $r$  dependent on the application).

#### 4.4. Spatial Correlation

In the presence of sudden events, the dynamics of a physical process can change rapidly. Often, detecting such events, such as a forest fire, a volcanic eruption, or a cardiac attack, is the very purpose of the WSN. However, the occurrence of the event may disrupt temporal correlations, giving rise to false anomalies. Nevertheless, different sensor nodes generally are affected by the event and produce measurements that are spatially correlated to the event source: as a consequence, the measurements of different sensors are correlated during the manifestation of the event [Boukerche 2009]. This phenomenon is known as spatial correlation.

Several techniques make use of spatial correlations by relating the measurements from different sensors in the same time interval—this is equivalent to fixing  $t$  in Equation (2) and letting the parameter  $i$  vary. The most widespread spatial correlation model is also the simplest: it assumes that all sensors would produce the same measurements in the absence of errors and noise; that is, they measure the same value, and we refer to this model as *spatially homogeneous* [Zhang et al. 2006; Ngai et al. 2006; Wu et al. 2007; Liu et al. 2007; Bettencourt et al. 2007]. In terms of the physical attribute model given in Equation (1),  $\varphi(s, t)$  is actually a function of time only. In this scenario, the sensors' measurements can be described by a Gaussian distribution. This

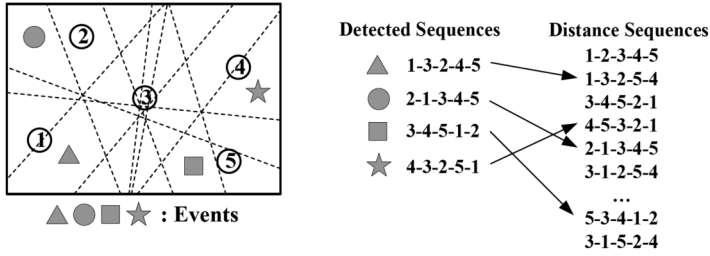


Fig. 2. Detection of measurements that do not comply with the monotonicity assumption, from Guo et al. [2009].

is because they are independent observations of random variables with a well-defined expected value and well-defined variance, and according to the *central limit theorem*, their values will be approximately normally distributed [Rice 2007]. Detecting sensors with abnormal readings becomes, then, a simple matter of detecting deviations from the spatial measurements' distribution, and the accuracy of the distribution estimation increases with the number of sensors.

The homogeneous model is suitable only for regions of space that are small enough and free of obstacles. However, when the deployment topology and characteristics of the physical phenomena violate the homogeneity assumption, the spatial propagation rules can still induce spatial correlations. In many applications, such propagation can be assumed monotonic [Guo et al. 2009]. This implies that the values of the physical attribute at a point in space should either increase or decrease as the distance from that point increases. For example, when monitoring for forest fires, the temperature decreases monotonically as the distance from the fire increases. To ascertain whether this property holds, Guo et al. [2009] divide the deployment space into sections, called faces. For each face, the authors construct a “distance sequence,” corresponding to the sequence of sensors ordered by the distance from that face. While sensing the phenomenon, the sensors' readings are sorted to generate the *estimated sequence*, which is then compared to all possible distance sequences, as shown in Figure 2. The sensors' measurements are consistent with the expectation if the estimated sequence corresponds exactly to one of the distance sequences. This condition is then relaxed to cope with noisy measurements that degrade the validity of the monotony assumption, but the main factor undermining its validity is the presence of multiple simultaneous events [Guo et al. 2009].

Instead of considering a strict assumption like the monotonicity of the measurements, it is possible to model correlations between the sensors' readings as a function of their spatial positions. An example of such a model is the *variogram*, defined as the variance of the difference between values of a physical phenomenon at two locations. In our notation, the variogram between two points  $s_1$  and  $s_2$  is defined as  $var(\varphi(s_1, t) - \varphi(s_2, t))$ . When the physical phenomenon is assumed to be isotropic, the variogram is expressed as a function of the distance only, and Zhang et al. [2012] have applied it to compute an expected measurement as a function of the measurements from other sensors. Note that in the presence of obstacles, the variogram not only is a function of the distance but also depends on the absolute positions.

Rather than considering distances between sensors, spatial correlation can be calculated as a function of the sensor values themselves. This choice caters for sensors at the same distance but subject to different noise or obstacles in space. However, it comes at the price of correlation updates when sensors are mobile. For example, Sharma et al. [2010] express a sensor's measurement as a linear combination of the measurements from the other sensors, extract the function's parameters, and derive expected

sensor readings. Dereszynski and Dietterich [2011], instead, derive expected readings by fitting the joint probability distribution of the measurements from  $N$  sensors, after assuming it is an  $N$ -variate Gaussian distribution. Note that this approach also implicitly assumes a linear model, as the covariance between two random variables captures linear dependencies (we have mentioned in Section 4.2 that this is true for the Pearson correlation coefficient, which is just a normalization of the covariance index).

Not infrequently, spatial correlation is used in conjunction with temporal correlation, since they capture different kinds of deviations. For example, Bettencourt et al. [2007] propose an outlier detection technique based on two kinds of differences: between a sensor's reading and its own previous reading (temporal correlation) and between readings of different sensors at the same time (spatial correlation). A distribution for both differences is used to check if data samples are statistically significant as related to the temporal domain as well as to the spatial domain.

#### 4.5. Attribute Correlation

In the same WSN, sensors observing different physical attributes such as light, vibrations, temperature, and so forth may coexist. Some of these attributes may be correlated because of the physical relationship between them, for example, temperature and relative humidity. Commonly, at every deployment location,  $s_i$  different sensors in charge of measuring different physical processes are connected to a single sensor node. As described by Equation (5), for a fixed point in space and time, we have a set of  $A$  physical attributes. We define attribute correlation as the correlation between them:

$$\varphi^a(s, t) \quad a \in 1, \dots, A. \quad (5)$$

We expect attribute correlations to also be observable in the measurements reported by the sensor nodes. Note, however, that attribute correlations between sensors belonging to the same node are not informative as an attacker who has compromised a node may tamper with all the measurements collected on that node. However, attribute-based expectations are very useful in conjunction with spatial correlations when spatial redundancy is limited. For example, body sensor networks for health care have limited redundancy since it is impractical to cover the patient with several sensors. We can then still exploit correlation among different physiological values (the attributes) measured by different sensor nodes.

An example in the health care domain is presented by Salem et al. [2013], who exploit spatial attribute correlations together with temporal correlations. Based on a Discrete Wavelet transform, they decompose the attribute signals into average and fluctuation signals. Abrupt temporal changes in the energy of the fluctuation signal are detected by a Hampel filter, which flags outlying attributes. This technique has been proposed for fault-tolerant event detection, based on the observation that multiple attributes are expected to be flagged simultaneously in the presence of an event, due to their attribute correlations. Then, if the minimum number of outlying attributes is not reached, the sensors reporting the outlying readings are considered faulty. However, in the context of malicious data injections, this technique would not prevent an attacker from deliberately injecting measurements that subvert the event detection.

#### 4.6. Overview of Techniques for Extracting Expected Data

In the previous sections, we have analyzed different types of correlations, the information they capture, and variations in the exploitation of the same correlation across the techniques proposed in literature. In Table I, we summarize this analysis.

Table I. Correlation Types

Correlation Type	Information Captured	Variations
Temporal	$\text{corr}(\varphi^a(s, t_1), \varphi^a(s, t_2))$	<ul style="list-style-type: none"> <li>— Time-series evolution model</li> <li>— Time memory (the maximum value of <math>W</math> for which the correlation is modeled)</li> </ul>
Spatial	$\text{corr}(\varphi^a(s_1, t), \varphi^a(s_2, t))$	<ul style="list-style-type: none"> <li>— Spatial model, e.g., homogeneous, monotonic, variogram, linear dependency</li> <li>— Correlation variational model, e.g., distance dependent, sensor dependent, fixed</li> <li>— Neighborhood selection criterion</li> </ul>
Attribute	$\text{corr}(\varphi^{a_1}(s, t), \varphi^{a_2}(s, t))$	<ul style="list-style-type: none"> <li>— Correlation extraction process, e.g., from physical laws, temporal/spatial analysis, etc.</li> </ul>

## 5. DETECTING DEVIATIONS FROM EXPECTED DATA

Expectations about the actual measurements can be used to calculate the deviation of the reported measurements from them. Both anomaly detection and trust management require an expectation, but they use different criteria to cope with abnormal data. Specifically, anomaly detection uses the expectation to discriminate between anomalous and normal data. Trust management instead uses a criterion to map the deviation from expected data to a trust value. Since the two techniques differ in how they interpret deviation, we will consider them separately in this section.

### 5.1. Anomaly Detection Techniques

Anomaly detection is a method to characterize data as normal or anomalous. In contrast to Rajasegarar et al. [2008], who consider outlier detection and anomaly detection as equivalent, we instead consider outlier detection as one of the techniques belonging to the anomaly detection category. The reason is that outlier detection identifies the samples that are unlikely to manifest. However, the measurements could be anomalous with respect to other criteria that cannot be reduced to the problem of finding outliers. Consider, for example, the case where a sensor is experiencing a *stuck at fault*; that is, it always produces the same measurement. An outlier detection technique applied on a subset of the last measurements from that sensor will detect no outlier. However, an anomaly still exists and could be detected by considering, for instance, the low variance in the measurements' distribution. To clarify this aspect, we present statistical tests for anomaly detection and highlight their differences with more traditional outlier detection techniques. Then we delve into techniques for outlier detection, which is still the most commonly adopted technique for anomaly detection.

*5.1.1. Statistical Tests.* Techniques based on statistical tests assume a probabilistic data distribution. Real data is then checked against this distribution to verify its compliance to it. Techniques based on statistical tests are more general than outlier detection because they check the compliance of both outliers and nonoutliers to the distribution, whereas outlier detection focuses on the *classification* of single data samples.

For example, Rezvani et al. [2013] use a technique based on statistical tests to detect malicious colluding nodes. They assume spatial homogeneity and model sensor measurements as a ground-truth value plus some noise. The ground truth is estimated as a weighted average of measurements, and the difference between the estimated value and each measurement is assumed to be normally distributed. This assumption is in keeping with the application of the *central limit theorem* [Rice 2007]—errors are assumed to be due to a large number of independent factors and thus to follow a normal distribution. Compliance with the normal distribution is then assessed with the

Kolmogorov-Smirnov test, which quantifies the distance between an empirical distribution (the errors distribution) and a reference distribution (the normal distribution).

*5.1.2. Outlier Detection.* Outlier detection methods consider as anomalous data that lies outside of the space where most data samples lie. This technique can identify malicious data injections reasonably effectively as long as maliciously injected values are a minority in the dataset and deviate significantly from the other data.

Historically, outlier detection has been proposed in WSN for different purposes, sometimes with opposing goals: in some cases, the techniques aim to filter out outliers, and in others the outliers represent the main interest. For example, outliers are filtered out to increase data accuracy [Janakiram et al. 2006] and for energy saving [Rajasegarar et al. 2007]. Applications where outliers are the main interest include fault detection [Paradis and Han 2007], event detection [Bahrepour et al. 2009; Zhang et al. 2012], and detection of malicious data. We describe next different approaches to the outlier detection problem independently of the application context, but we focus on those techniques that can be applied to detecting malicious data injections.

*Nearest-Neighbor-Based Outlier Detection.* In nearest-neighbor-based outlier detection, an outlier is a data sample with a narrow neighborhood, where a neighborhood comprises the data samples within a certain distance. Most nearest-neighbor-based techniques in WSNs are inspired from the well-known LOCI method [Papadimitriou et al. 2003], which calculates, for every sample, the number of neighbors in a data space characterized by the radius  $\alpha r$ , where  $\alpha$  is a parameter used to reduce computational complexity. The relative difference with the average number of neighbors, that is, the samples within a radius  $r$  in the data space, constitutes the *Multigranularity Deviation Factor* (MDEF). The MDEF is compared to a threshold equal to 3 times the MDEF standard deviation to ensure that less than 1% of values are above the threshold when the distances between data samples follow a Gaussian distribution (the percentage increases up to 10% for other distributions). Note that this method is applicable to malicious data injections by considering the sensors' measurements as the data samples. However, the research community seems to have lost interest somewhat in approaches based on nearest neighbor since they have large computational overheads due to the calculation of the neighbors for each new data sample.

*Clustering-Based Outlier Detection.* Clustering is another technique often used for outlier detection. Here the outliers are the elements distant from the others, after organizing close elements into clusters. For example, Rajasegarar et al. [2006] identify a cluster as anomalous if its distance to other clusters is more than one standard deviation of the distance of the cluster elements from the mean.

*PCA-Based Outlier Detection.* PCA [Marsland 2009] is a common data analysis technique that has also been applied to find outliers [Chatzigiannakis and Papavassiliou 2007]. PCA is based on a projection of the  $k$ -dimensional data space onto another  $k$ -dimensional data space, where the variables describing the data samples are linearly uncorrelated. This transformation is defined in such a way that the projected variables are sorted with descending variance. The first  $p$  out of  $k$  variables are defined as the *principal components* and can be projected back to the original data space to obtain a prediction vector  $y_{norm}$  [Jackson and Mudholkar 1979], also referred to as *normal data* [Chatzigiannakis and Papavassiliou 2007]. The difference between original and normal data constitutes the *residual vector*  $y_{res}$ . Residual vectors that are large in magnitude (i.e., when the squared prediction error  $SPE = \|y_{res}\|^2$  of the residual vector is greater than a threshold) are interpreted as deviations from the predicted (normal) vector and considered as outliers [Chatzigiannakis and Papavassiliou 2007]. PCA can be applied to  $k$ -dimensional datasets, for example, made up of the measurements' time series of  $k$

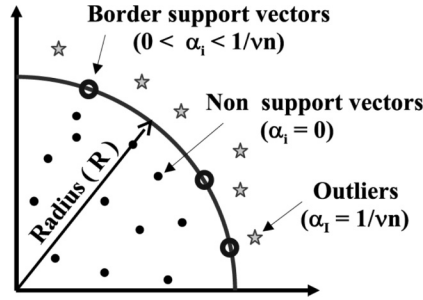


Fig. 3. One-class quarter-sphere support vector machine, from Rajasegarar et al. [2007].

sensors [Chatzigiannakis and Papavassiliou 2007]. In this case,  $y_{res}$  reflects changes in spatial correlation, but the same idea can also be applied to the temporal or attribute domains.

*Classification-Based Outlier Detection.* Traditional classification techniques learn how to recognize samples from different classes. Anomaly detection considers two classes: anomalous and normal; however, anomalous data samples are rarely observable compared to the normal ones. Therefore, classification for anomaly detection is generally reduced to a one-class classification problem, based on the observation of normal samples only.

Normal and anomalous samples can be viewed as points within two different regions of the data space. Finding the boundary that separates the two regions may be infeasible, because the regions overlap and, even when a boundary exists, it may have a complex shape. Support Vector Machine (SVM) is a classification technique that can overcome this limitation by projecting the data samples into a higher-dimensional space. In the projected data space, a boundary that separates normal from anomalous points may exist even if it does not exist in the original space, or may have a simpler shape. For example, the normal samples could be contained within a sphere in the projected data space. When the data space contains only positive values, this problem reduces to a special type of SVM called *one-class quarter-sphere SVM* [Laskov et al. 2004], which is represented in Figure 3. With this approach, the classification problem reduces to finding the sphere’s radius. Depending on how the WSN dataset is given in input to quarter-sphere SVM, the classification can be made across its time domain [Rajasegarar et al. 2007], attribute domain, or both [Shahid et al. 2012].

Bayesian networks have also been applied in WSNs to detect outliers with a classification-based approach. A Bayesian network defines the relations of conditional independence between random variables through a network graph. In WSNs, the random variables can be different values in space and time of the physical attributes.

An example of application of Bayesian networks to WSNs is given by Dereszynski and Dietterich [2011]. The physical attribute  $\varphi(s_i, t_k)$  is modeled as a random variable that depends on  $\varphi(s_i, t_{k-1})$  (first-order Markov relationship) and on values at different locations  $\varphi(s_{j \neq i}, t_k)$ . The aim is to find the state of a sensor, modeled by a random variable with two possible values: *working* and *broken*. The posterior probability of the measurements, which depends on both the physical attribute and the sensor state variable, is maximized with respect to the state variables to identify faulty nodes. Dereszynski and Dietterich [2011] evaluated their approach assuming that faulty sensors have a high increase in their measurements’ variance (by  $10^5$ ), motivated by the observation that the measurements of faulty sensors appear more noisy. Though reasonable in the case of faults, this assumption usually does not hold for data injections, where an

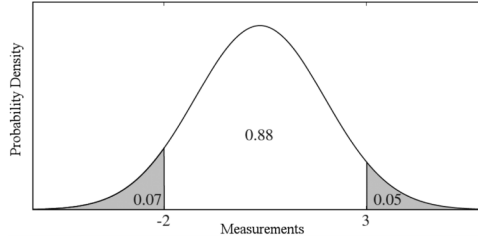


Fig. 4. Statistical characterization of the sensed data for outlier detection, from Bettencourt et al. [2007].

attacker can choose the measurements' distribution arbitrarily and wishes in most cases to remain undetected.

*Statistical Outlier Detection.* Statistical outlier detection identifies outlying data samples through statistical characterization of the tail of the samples' probability distribution, as shown in Figure 4.

Note that this approach differs from anomaly detection based on statistical tests, as it does not test the samples' compliance to their expected distribution but only identifies the outliers that lie on the tails of the distribution. For example, outliers can be defined as samples far from the mean. Ngai et al. [2006] have applied this idea to measurements from different sensors, thus exploiting spatial correlation. The *spatial* sample mean  $\hat{\mu}_S$  of measurements from  $N$  different sensors is defined as

$$\hat{\mu}_S = \frac{1}{N} \sum_{j=1}^N S_j(t). \quad (6)$$

Ngai et al. [2006] use it to evaluate the deviation of sensor  $j$  from the spatial mean, compared to the magnitude of the mean itself, with the metric  $f(j, t) = \sqrt{\frac{(S_j(t) - \hat{\mu}_S)^2}{\hat{\mu}_S}}$ .

Similarly, Tanachaiwiwat and Helmy [2005] use the metric  $t^* = \frac{S_i(t) - (\mu_{T_i} \pm \delta)}{S_{T_i} / \sqrt{W}}$ , where  $\mu_{T_i}$  and  $S_{T_i}$  are respectively  $i$ 's temporal mean and sample standard deviation in a window of size  $W$ , and  $\delta$  is an already known variation between sensor  $i$  and  $j$  due to the observed phenomenon's spatial propagation. Considering the model in Section 4.1, a generic sensor  $j$  calculates its *temporal* sample mean in the  $W$ -wide time window  $[t_{K-W+1}, t_K]$  as

$$\hat{\mu}_{T_j} = \frac{1}{W} \sum_{n=0}^{W-1} S_j(t_{K-n}). \quad (7)$$

The *temporal* standard deviation is instead calculated as

$$S_T = \sqrt{\frac{1}{W-1} \sum_{n=0}^{W-1} (S_j(t_{K-n}) - \hat{\mu}_{T_j})^2}. \quad (8)$$

The value of  $t^*$  is then compared with a threshold that is set to 3 since, in normally distributed data, this accounts for approximately 99.7% of the population (the percentage decreases to 90% for other distributions).

In some cases, the median is preferred to the mean, since the former has the advantage of being insensitive to outliers. Indeed, a problem in outlier detection is how to find the general (nonoutlying) trend from data affected by outliers. The mean is

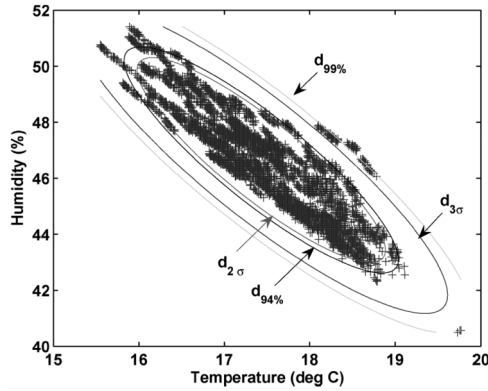


Fig. 5. Statistical distribution in the attribute space made up by temperature and humidity. Points with Mahalanobis distance greater than  $d$  are treated as outliers, from Rajasegarar et al. [2009].

sensitive to outliers, since it is proportional to the magnitude of each operand. The median takes instead one element to represent all the others. Wu et al. [2007] use the median operator to aggregate sensors' measurements in a neighborhood. We can refer to it as a *spatial median*. If we order the  $N$  sensors' measurement at time  $t$  such that  $S_1(t) \leq S_2(t) \leq \dots \leq S_N(t)$ , the median in the spatial domain is calculated as:

$$\tilde{\mu}_S = \begin{cases} S_{(N+1)/2}(t) & \text{if } N \text{ is odd} \\ S_{N/2}(t) & \text{if } N \text{ is even.} \end{cases} \quad (9)$$

After calculating the difference between the median and each value, there are two possibilities: comparing each difference to the measurements' magnitude or comparing it to the general distribution of the differences. Yang et al. [2006] and Wu et al. [2007] detect outliers in the differences, assuming they are normally distributed. Instead of relying on the assumption of a Gaussian distribution, the probability distribution can also be estimated from the data [Bettencourt et al. 2007].

When sensing multiple physical attributes, the distribution of the measurements across all attributes can be considered, rather than a separate distribution for each one. This approach can potentially detect outliers that a separate approach would fail to detect. Liu et al. [2007] combine different attributes using the Mahalanobis distance, which is based on the interattribute correlation and defines how the data is statistically distributed in the attribute space. This scheme is shown in Figure 5.

## 5.2. Trust-Management-Based Techniques

Trust management considers the trustworthiness between two classes of entities: a trustor and a trustee. The trustor assigns each trustee a trustworthiness value, based on how much the trustee's behavior matches an expectation. Trustworthiness values are usually in the range  $[0, 1]$ , decreasing when the trustee exhibits deviations from the expected behavior and increasing when the trustee's behavior matches the expectation.

Trust management can be usefully applied in WSNs to reduce the influence of the compromised sensor nodes that inject malicious data. Indeed, if the expected behavior accurately characterizes genuine nodes, compromised nodes would be assigned a low trustworthiness value when deviating from it. Since trust values are a continuous metric defined inside an interval, there is no direct classification of compromised and genuine nodes. Instead, the trust values are used to apply a penalization proportional to the confidence that the sensor is compromised. Note that the influence of the compromised nodes becomes negligible only when the confidence is sufficiently high. Filtering



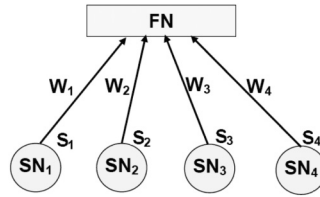


Fig. 6. Trust-weighted aggregation for event detection. FN is a forwarding node, which collects reports from the sensor nodes SN, from Atakli et al. [2008].

all the sensors with a trustworthiness under a given threshold [Sun et al. 2012] could help mitigate this drawback but requires a method to set the appropriate threshold.

**5.2.1. Event-Based Techniques.** Trust-management for sensed data was originally introduced as a complement to network-level trust, that is, how much nodes can be trusted to perform correctly network-level tasks [Ganerival et al. 2003; Raya et al. 2008; Momani et al. 2008] such as communicating routes, participating in the route discovery process, routing incoming packets, and so forth. The behavior with respect to each of these tasks can be of two kinds: cooperative and uncooperative.

The first examples of trust management for sensed data use a similar binary evaluation to build the trustworthiness, defined with respect to an event detection process. Initially, a *decision logic* establishes the presence of the event by combining the sensed data and the trust values. Then, the sensed data is compared to the final decision to measure the sensor’s cooperativeness and update the trust values. This criterion is based on the assumption that nearby sensors are expected to agree about the event presence, which is a form of spatial correlation (see Section 4.4).

One of the first techniques to adopt this approach is described in Atakli et al. [2008]. As shown in Figure 6, the reading of a generic sensor  $i$ ,  $S_i(t)$ , which can take the values 0 and 1 (absence/presence of an event), are relayed to a *forwarding node*. This node computes  $\sum_{n=1}^N W_n S_n(t)$ , where  $W_{n:n \in 1 \dots N}$  denote the trust weights.

The result is used to decide about the ground truth  $E$ . Afterward, weights are updated with the following rule:

$$W_n = \begin{cases} W_n - \theta r_n, & \text{if } S_n(t) \neq E \\ W_n, & \text{otherwise,} \end{cases} \quad (10)$$

where  $r_n$  is the ratio of sensors giving different outputs over the total number of sensors and  $\theta$  is a penalty weight that determines a tradeoff between the detection time and accuracy. In summary, the trustworthiness values, which coincide with the weights, are calculated based on the measurements’ consistency with the aggregated value. The latter is considered more reliable than the single readings, since sensors that exhibited inconsistent (e.g., malicious) readings in the past contribute less to the aggregation process. Finally, malicious nodes are detected by comparing the weights to a threshold, which the authors heuristically set to 0.4. Note that the algorithm is vulnerable to the *on-off attack*: a node that performs well for a time period acquires high trustworthiness and then suddenly starts malfunctioning [Sun et al. 2006].

To counteract the on-off attack, Oh et al. [2012] and Lim and Choi [2013] propose to penalize  $S_n(t) \neq E$  by a quantity  $\alpha$  and reward  $S_n(t) = E$  by a quantity  $\beta$  with  $\beta < \alpha$ . As  $\frac{\alpha}{\beta}$  grows bigger, faulty and malicious nodes are filtered out faster. However, sensors with transient faults are also filtered out, even though they may report correct measurements later on. To avoid this, the ratio  $\frac{\alpha}{\beta}$  needs to also consider the probability of transient faults and their duration distribution. Therefore, this operation just reduces

the frequency with which an attacker can switch between “good” and “bad” behavior in an on-off attack.

When the sum of all trust weights is equal to 1, the weighted sum of sensors’ reading corresponds to a weighted mean. As described in the previous section, the mean has the drawback of being directly proportional to extreme readings. So in trust-based aggregation as well, the median could be used as a more robust aggregation operator. A trust-weighted median has been applied by Wang et al. [2010] in the context of acoustic target localization, where the median allows one to filter out faulty measurements. The advantages of using the weighted median increase when an element with high weight has an extreme value. Indeed, while the weighted mean would be biased toward that value, the weighted median would still filter it out, if the other values are not extreme and the sum of their weights is bigger than the weight of the extreme value. This property reduces the efficacy of an on-off attack.

Another aspect to take into account is the uncertainty in the event’s presence. Raya et al. [2008] deal with this aspect by using a decision logic based on Dempster-Shafer Theory (DST), which expresses the belief about the event presence as a combination of individual beliefs from sensor nodes. DST combines the sensors’ information supporting the event with the information not refuting the event (the uncertainty margin that may comply with the event presence).

*5.2.2. Anomaly-Based Techniques.* Rather than analyzing the compliance with the output of an event decision logic, other trust management techniques look for anomalous behaviors with techniques similar to anomaly detection ones.

In fact, the output of anomaly detection itself can be used to define a cooperative/uncooperative behavior [Ganeriwal et al. 2003], but a more flexible approach that does not restrict the observations to a binary value is to update trust values based on an anomaly score. An example is given by Bankovic et al. [2010], using self-organizing maps (SOMs). The SOM is a clustering and data representation technique that maps the data space to a discrete 2D neuron lattice. Bankovic et al. [2010] build two SOM lattices: one in the temporal domain and another in the spatial domain. The trust values are assigned based on two anomaly scores: the distance between the measurement and the SOM neuron, and the distance between the neuron to which the measurement has been assigned and other SOM neurons. The main disadvantage of this algorithm is its computational time. For better accuracy, SOMs require many neurons, but the computational time increases noticeably [McHugh 2000].

Another example is given by Zhang et al. [2006], who use a statistical-test approach (see Section 5.1.1) to assign reputation values to the sensors. The measurements gathered in time are assumed to approximately follow a normal distribution. The normal and actual measurements’ distribution are compared with the Kullback-Leibler divergence  $D_n$ , which evaluates the information lost when a probability distribution is used in lieu of another. The divergence is then used to update the trust values, with the following expression:

$$W_n = \frac{1}{1 + \sqrt{D_n}}. \quad (11)$$

*5.2.3. Using Second-Hand Information.* In the trust management schemes previously analyzed, each sensor’s trust values are computed and updated by the device with the trustor role, typically a forwarding node. However, when the trustor is not in the transmission range of its trustee  $i$ , it may rely on information from its neighbors  $N_i$  to calculate its trustworthiness. Bao et al. [2012] deal with this problem by introducing two different trust update criteria:

$$T_{ij}(t) = \begin{cases} (1 - \alpha)T_{ij}(t - \delta t) + \alpha T_{ij}(t) & \text{if } j \in N_i \\ \text{avg}_{k \in N_i} \{ (1 - \gamma)T_{kj}(t - \delta t) + \gamma T_{kj}(t) \} & \text{otherwise.} \end{cases} \quad (12)$$

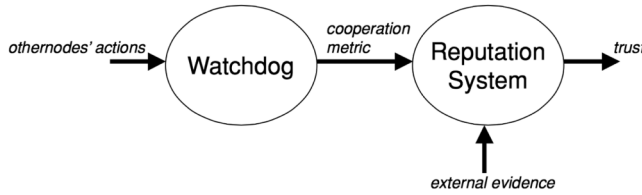


Fig. 7. Combination of direct information and recommendations, from Ganeriwal et al. [2003].

The calculations of the second case represent node  $j$ 's *recommendation*, that is, the trustworthiness extracted from relayed information. Eventually, recommendations depend on trustworthiness from the viewpoint of direct neighbors. However, such trustworthiness can be manipulated by malicious nodes to bad-mouth or good-mouth other nodes. Bao et al. [2012] mitigate this problem by controlling the impact of recommendations through parameter  $\gamma$ , set to  $\frac{\beta T_{ik}(t)}{1 + \beta T_{ik}(t)}$ . Thus, if a sensor has little trust compared to the parameter  $\beta$ , the contribution of its recommendation will be small. However, sensors conducting an on-off attack can give false recommendations for a short while and then behave correctly again without being detected.

Even when direct information is available, recommendations can be used as second-hand information and combined with direct information to obtain a *reputation*. Second-hand information speeds up the convergence of trust values but adds network traffic overhead and introduces new problems, such as the weighting criterion for recommendations and the recommendation exchange frequency [Huang et al. 2006]. Ganeriwal et al. [2003] follow this approach and treat reputation as a probability distribution, updated as a combination of direct and indirect reputation. Direct reputation is updated based on a watchdog module, while indirect reputation is updated with recommendations, that is, reputation from other nodes. The framework's scheme is shown in Figure 7. Note that such definition of reputation introduces a loop: indirect reputations come from reputations given by other sensors, which in turn depend on indirect reputations. To avoid the information loop, the recommendations need to be taken only from direct observers.

Modeling the reputation as a single value does not consider the uncertainty that a sensor has in trusting another sensor. This information is particularly useful with recommendations, as recommendations from sensors with high uncertainty should contribute less. To consider uncertainty, the reputation can be modeled with a probability distribution, whose choice is dictated mainly from the trust evaluation and update criteria. For example, Ganeriwal et al. [2003] use the *beta distribution* since it is the posterior distribution when the binary interactions between nodes are modeled with a binomial distribution. Momani et al. [2008] apply a normal distribution to model the differences between the measurements of two sensors (spatial homogeneity is assumed; see Section 4.4).

## 6. DIAGNOSIS AND CHARACTERIZATION OF MALICIOUS DATA INJECTIONS

Detecting the deviation of the measurements from the expected behavior is usually not sufficient to infer the presence of a data injection attack. In the case of outlier detection, for example, we have seen that measurements are only classified as outlying or nonoutlying, but malicious data injection is only one of the possible causes for outlying data. In general, regardless of the technique that detects the deviation from an expected behavior, the cause for that deviation needs to be found. We refer to this task as *diagnosis*. Generally, it is not a trivial task, because different causes such as faults or genuine events may have similar effects.

Additionally, even when the presence of an attack can be ascertained with confidence, further information is needed to determine the course of action to be taken. For example, there is the need to know the attack's effects and the system area (nodes) affected by the attack. We refer to this other task as the *characterization* of the attack.

In the following, we analyze the state of the art for diagnosis and characterization of malicious data injections in WSNs.

### 6.1. Diagnosis

Diagnosis of malicious data injections in WSNs consists of distinguishing them from two main phenomena that can produce similar deviations from expected behavior: faults and events of interest. Faults represent generic unintentional errors introduced, for example, by obstacles in the environment, sensors' battery depletion, pollution, and fouling. Events of interest represent environmental conditions that seldom manifest but are interesting as they can reveal an alarm scenario, for example, heart attacks, fires, and volcanic eruptions.

Information about the cause of an anomaly or of an untrustworthy sensor can be precious. With fine-grained knowledge about the nature of the problem, an appropriate response can be initiated to address it. Unfortunately, in the papers analyzed so far, an exhaustive diagnosis phase is still lacking. Most of the attention has focused on diagnosing events as opposed to faults. The general assumption used to distinguish between them is that faults are likely to be stochastically unrelated, while event measurements are likely to be spatially correlated [Luo et al. 2006; Shahid et al. 2012]. Note that this assumption excludes common-mode failures from the analyses. Based on this assumption, after detecting deviations from expected data with temporal [Bettencourt et al. 2007; Shahid et al. 2012] or attribute [Shahid et al. 2012] correlations, it is possible to diagnose whether the deviation was caused by a fault or an event by exploiting spatial correlation. When there is a consensus among a set of sensors about the presence of an event, discording sensors are considered faulty [Luo et al. 2006; Shahid et al. 2012; Bettencourt et al. 2007]. Similarly, some sensed attributes (e.g., human vital signs, such as glucose level, blood pressure, etc.) can be assumed heavily correlated in the absence of faults, which instead disrupt attribute correlations. Then, if we further assume that events would cause a minimum number of outlying attributes, faults can be identified when the minimum is not reached [Salem et al. 2013].

Fewer advances have been made toward diagnosing malicious interference as opposed to faults and events. We summarize them in the following sections.

*6.1.1. Distinguishing Malicious Interference from Events.* In the literature, malicious interference is distinguished from events through an agreement-based strategy [Liu et al. 2007; Atakli et al. 2008; Wang et al. 2010; Oh et al. 2012; Lim and Choi 2013; Sun et al. 2013]; that is, the sensor's information is first used to decide about the presence of an event and then sensors that did not support the final decision are identified as malicious. This approach is based on the assumption that sensors are sufficiently spatially correlated to correctly detect events. However, multiple compromised nodes can also collude in the attack to keep the spatial correlations consistent between themselves. This complicates discriminating between genuine events and malicious data injections and allows an attacker to fabricate false events or to mask genuine ones. This aspect is discussed in more detail in Section 6.2.

*6.1.2. Distinguishing Malicious Interference from Faults.* Criteria to distinguish malicious data injections from faults are less remarked on. Two main approaches can be identified: delegating the diagnosis to intrusion detection techniques and leveraging fault statistics.

*Intrusion Detection.* One of the main challenges in detecting attacks with anomaly-based techniques is that such techniques abstract the means through which an attack is conducted. This choice comes from their objective to detect new attacks with unknown patterns, as opposed to intrusion detection techniques, which are based on recognizing known attack signatures. The framework proposed by Ngai et al. [2006] is a tradeoff between an anomaly detection technique and an intrusion detection system, since the detection is carried out through anomaly detection achieving a high detection rate, while the diagnosis is carried out with intrusion detection. Clearly this approach provides diagnosis only for known attacks and cannot distinguish between an unknown attack and a fault.

*Fault Statistics.* The statistical characterization of faults can also be used to distinguish them from malicious interference. Oh et al. [2012] and Lim and Choi [2013] use the expected frequency of transient faults to avoid excluding from the system sensors subject to transient faults. Indeed, their trust management algorithm allows such sensors to recover trustworthiness by allowing temporary misbehavior. Only sensors misbehaving with higher frequency, including malicious sensors and sensors with permanent faults, will then be excluded.

## 6.2. Characterization

If *detection* and *diagnosis* of malicious data injections answer the question “Is there an attack?”, *characterisation* answers questions such as “Which are the compromised sensors?” and “How is the attack performed?” The difference is perhaps more evident in event detection tasks. For example, *after* detecting the presence of an event, the event’s spatial boundary can be characterized using the methodology proposed in Wu et al. [2007], which finds the areas where the difference between the measurements from different sensors is high, indicating a discontinuity introduced by the event boundary. In this case, characterization is triggered by detection but is a separate task.

*6.2.1. Collusion and Its Effects.* In malicious data injections, detection, diagnosis, and characterization are often addressed simultaneously, since the information characterizing the attack can be precious to improve the detection. In particular, when multiple sensors have been compromised and *collude* in the attack, they act in concert to change the measurements while evading, if possible, any anomaly detection applied. Therefore, identifying which sensors are more likely genuine and which sensors are more likely compromised becomes an integral part of detecting the attack itself.

In *collusion attacks*, compromised sensors follow a joint strategy that reduces the advantages of spatial correlation, since the compromised nodes cooperate to form credible spatially correlated data [Tanachaiwiwat and Helmy 2005]. In the presence of collusion, diagnosis is also significantly more complex. Tanachaiwiwat and Helmy [2005] point out that when a genuine outlier (e.g., related to an event) occurs, extreme readings from the colluding nodes could be hidden. The problem becomes increasingly difficult as the percentage of (colluding) compromised sensors increases. Ultimately, when the number of colluding sensors increases to the point of exceeding genuine sensors, the attack may still be detectable, but it may be impossible to distinguish which nodes are genuine and which nodes are compromised. Tanachaiwiwat and Helmy [2005] evaluate their anomaly detection algorithm against colluding nodes and find that performance noticeably decreases when more than 30% of the nodes are colluding. A similar result is reported by Chatzigiannakis and Papavassiliou [2007].

Bertino et al. [2014] describe a new attack scenario applicable when the trustworthiness is calculated through an *iterative filtering* algorithm. While in generic (non-iterative) trust evaluation techniques trust weights are updated based on data from the current time instant and the weights calculated at the previous time instant, in

iterative filtering, the weights are iteratively updated with data of the same time instant until a convergence criterion is satisfied. In this context, the authors introduce a new attack scenario where all colluding nodes but one produce noticeable deviations in their readings. The remaining compromised node reports, instead, values close to the aggregated value of all the readings (including malicious ones). Eventually, this node acquires a high trust value, while all the others acquire low trust values. The aggregated value, in turn, quickly converges to a value far from that of the genuine nodes. The authors show that this attack is successful when the sensors are assigned equal initial trustworthiness. They therefore propose to calculate the initial trustworthiness as a function that decreases as the error variance increases. The error is defined as the distance from an estimated physical attribute value  $\varphi(t)$  and is the same for all the sensors.

Rezvani et al. [2013] proposed another technique that detects collusion rather than counteracting it. This technique is based on the assumption that deviations from the aggregated values are normally distributed for genuine nodes. This assumption comes from the observation that the deviations of noncompromised nodes, even if large, come from a large number of independent factors, and thus must roughly have a Gaussian distribution. For colluding nodes, instead, they assume that this condition does not hold. Then, by running the Kolmogorov-Smirnov test to check compliance to the normal distribution, they discriminate colluding nodes from genuine nodes.

In summary, while many studies propose new anomaly detection algorithms to cater to a broad range of scenarios, comparatively fewer address specifically malicious data injections in a way that can cater to more sophisticated attacks involving collusion between sensors. Such scenarios will need to be explored further in the future.

*6.2.2. Characterization Architectures: Centralized Versus Distributed.* To detect, diagnose, and characterize the nodes injecting malicious measurements, different architectures can be employed with different degrees of distribution. We discuss the properties of different solutions next.

In WSNs, there is always at least one entity that eventually collects the measurements for the analyses, decisions, and actions that the system needs to carry out: the basestation. The basestation is usually assumed free of compromise and therefore can be used to characterize the compromised nodes. In this case, we have a centralized architecture such as in Chatzigiannakis and Papavassiliou [2007], Atakli et al. [2008], Oh et al. [2012], Lim and Choi [2013], and Rezvani et al. [2013].

Even when the basestation is the only trusted entity in the network, distributed characterization is possible. Indeed, as proposed in Bao et al. [2012], the sensor nodes can be assessed in a hierarchical structure, where each node assesses the trustworthiness of nodes below it in the hierarchy. The basestation thus trusts nodes when a chain of trust can be established from that node to the basestation.

When the distribution principle is taken to the extreme, each node acts as a watchdog for all its neighbors and reports alerts to the basestation (or the next node in the hierarchy) [Ganeriwat et al. 2003; Tanachaiwiwat and Helmy 2005; Liu et al. 2007]. After all the reports are collected, a decision is taken based on algorithms such as majority voting [Hinds 2009]. The drawbacks of this approach are that it lacks global knowledge and for this reason is less robust to collusion attacks, and that it introduces significant network overhead given by the watchdog reports. Tanachaiwiwat and Helmy [2005] propose to overcome these problems by deploying multiple reliable tamper-resistant sensor nodes that probe suspicious nodes. This solution, however, requires additional expensive hardware, which undermines the cost advantages of WSNs.

## 7. DISCUSSION

In the previous sections, we have seen how different techniques can be applied to detect malicious data injections, how they leverage measurements' correlations, and the assumptions on which such correlations are based. We have examined the different detection techniques and how they find deviations from the expected behavior. We have highlighted the importance of distinguishing between different sources of deviations and presented the main directions of work toward this objective so far.

We now combine these analyses by building direct comparison tables, which summarize their main characteristics. A summary of the results reported by each of the techniques mentioned is provided in the following section.

### 7.1. Comparison of Approaches

We divide our comparison of the approaches analyzed so far into Tables II and III, containing the anomaly detection and trust management techniques, respectively. The content of the columns from left to right is as follows: technique name and reference; correlation used to define expected data; assumptions about the spatial model if any; detection criterion used; possible sources of anomalies (as mentioned in their paper); and for which of them diagnosis criteria are given—for example, {Event},{Malicious or Faulty} means that the authors give a criterion to discern between anomalies arising from events and from malicious or faulty sensors.

We observe that spatial correlation is most often exploited, and this under the frequent assumption of a homogeneous space. The situation is particularly evident for papers considering the presence of malicious data injections and probably a consequence of the fact that, generally, only a minor subset of sensors is assumed to be compromised. Therefore, in the spatial domain, there is always a significant set of genuine measurements that can be exploited to detect the malicious ones.

Assuming spatial homogeneity makes the calculations significantly simpler, since the sensors are considered to measure the same value. However, it also significantly restricts the applicability of the techniques in real cases. When the physical phenomenon is observed with low precision, for example, overall temperature across a large open-space area, this assumption is still valid if the spatial variations are absorbed by the error term in Equation (2). However, this allows an attacker to introduce malicious data that is within the error bounds yet still deviates significantly from the real values. While this assumption is generally appropriate in small areas, small areas also typically include fewer sensors that have a higher risk of an attacker compromising them all.

When multiple types of correlation are considered, temporal correlations are generally exploited along spatial ones. Use of attribute correlations is rather infrequent, probably because understanding them requires knowledge about their physical significance and this is application specific. The tables highlight even more the lack of diagnosis and characterization (see Section 6.1). Few papers consider specifically malicious injections with collusion, and even fewer papers deal with the problem of distinguishing them from other causes of deviations. While distinguishing events from faults is the diagnosis more frequently considered, distinguishing attacks from faults is undoubtedly more challenging and still rather rare.

### 7.2. Comparing Reported Evaluation Results

In the previous sections, we have considered techniques that could be applied to the problem of detecting, diagnosing, and characterizing malicious data injections. For those techniques that focus specifically on malicious data injections, we now present

Table II. Anomaly Detection Techniques

Work	Correlation Exploited	Spatial Model	Detection Method	Classes Considered	Interclass Discrimination
EKF, CUSUM GLR [Sun et al. 2013]	Temporal	None	Change in the distribution of error from estimate	Event, Malicious, Faulty	{Event}, {Malicious or faulty}
MGDD [Subramaniam et al. 2006]	Temporal	None	Measurement probability	Event, Fault	None
Ngai et al. [2006]	Spatial	Homogeneous	Difference with neighbors	Suspicious of Sinkhole Attack	None
Wu et al. [2007]	Spatial	Homogeneous	Difference with neighbors	Event	None
FIND [Guo et al. 2009]	Spatial	Monotonic WRT event source	Spatial monotonicity disruptions	Fault	None
Salem et al. [2013]	Attribute-temporal	None	Energy of fluctuations	Event, Fault	{Event} {Faulty}
STIOD [Zhang et al. 2012]	Spatio-temporal	Variogram	Difference with estimate	Event, Error	{Event} {Error}
MAP+HBST [Ni and Pottie 2012]	Spatio-temporal	Linear spatial trend	Difference with estimate	Fault	None
Liu et al. [2007]	Spatial	Homogeneous	Difference with neighbors	Malicious, Event	{Malicious}, {Event}
ART [Tanachaiwiwat and Helmy 2005]	Spatial	Homogeneous	Difference with neighbors	Compromised, Uncalibrated, Sybil	{Compromised or Faulty}, {Uncalibrated}, {Sybil}
Rajasegarar et al. [2007]	Spatio-temporal	Homogeneous	Values outside a quarter-sphere	None	None
STA-QS-SVM [Shahid et al. 2012]	Spatio-temporal and spatio-attribute	Homogeneous	Values outside a quarter-sphere	None	None
Chatzigiannakis and Papavassiliou [2007]	Spatial	High Pearson correlation	Changes in correlation	Fault, Malicious	{Point failure or malicious node}, {Group failure or Collusion}
Bettencourt et al. [2007]	Spatio-temporal	Homogeneous	Distribution of temporal and spatial differences	Event, Fault	{Event}, {Point failure}
Handschin et al. [1975]	Spatial	Linear combination of state variables	Difference with estimate	Fault	None
Robust IF [Rezvani et al. 2013]	Spatial	Homogeneous	Distribution of distance from estimation	Fault, Malicious	None

the experimental evaluation setup used by the authors and compare the reported results. None of these techniques has been tested on real attack scenarios. This is not surprising as finding real attack data in existing WSN deployments is difficult. In fact, two approaches have been broadly adopted to evaluate the algorithms for detection of



Table III. Trust-Based Detection Techniques

Work	Correlation Exploited	Spatial Model	Detection Method	Classes Considered	Interclass Discrimination
Zhang et al. [2006]	Spatio-temporal	Homogeneous	Distance from mean of top-trust sensors	Malicious	None
WTE [Atakli et al. 2008]	Spatial	Homogeneous	Trust under a threshold	Malicious	None
Momani et al. [2008]	Spatial	Homogeneous	Trust under a threshold	Faulty, Malicious	None
Wang et al. [2010]	Spatial	Homogeneous	Difference with aggregated value	Faulty, Event	{Faulty}, {Event}
Bankovic et al. [2010]	Spatio-temporal	Heterogeneous	Difference with learned pattern	Malicious	None
Trust-based IDS [Bao et al. 2012]	Spatial	Homogeneous	Trust under a threshold	Malicious, Event	{Malicious}, {Event}
DWE [Oh et al. 2012]	Spatial	Homogeneous	Trust under a threshold	Malicious, Permanent Fault, Transient Fault, Event	{Malicious or Permanent Fault}, {Event}
Dual threshold [Lim and Choi 2013]	Spatial	Homogeneous	Trust under a threshold	Malicious, Permanent Fault, Transient Fault, Event	{Malicious or Permanent Fault}, {Event}

malicious data injections: *simulation* [Sun et al. 2013; Liu et al. 2007; Rezvani et al. 2013; Atakli et al. 2008; Bankovic et al. 2010; Oh et al. 2012; Bao et al. 2012; Lim and Choi 2013] and *injection of attacks* in real datasets [Tanachaiwiwat and Helmy 2005; Chatzigiannakis and Papavassiliou 2007].

Table IV summarizes all the results achieved, together with all the relevant simulation parameters. The last three columns express the false-positive rate (FPR) when the detection rate (DR) is respectively 0.90, 0.95, and 0.99. DR is, by definition, the number of attack instances that are correctly detected divided by the total number of attack instances. FPR is, by definition, the number of times normal data instances are misclassified as attacks divided by the total number of normal data instances. The relationship between DR and FPR is known as the Receiver Operating Characteristic (ROC). Column 2 reports information about the size of the dataset used in the experiments. Column 3 reports the percentage of either malicious nodes or malicious measurements. Column 4 reports the input size for the algorithm; for example, in an experiment with 100 nodes, where the nodes are clustered in groups of 10 and the algorithm is run on clusters, the algorithm input size is 10.

Generally, in each paper, the tests are conducted in scenarios with different assumptions. For instance, Liu et al. [2007] generate data with a normal distribution for normal sensors and another normal distribution for malicious sensors. The results are excellent but depend a lot on the difference between the two distributions. Another important assumption, which has a noticeable impact on the results, is the spatial model. As pointed out in Section 4.4, most papers assume that the sensors' readings are homogeneous in the space; in other words, the measurements are expected to be

Table IV. Detection Performances, Independent Attacks

Work	Dataset Size	Dataset Malicious Percentage	Input Size for Each Algorithm Execution	FPR for DR=0.90	FPR for DR=0.95	FPR for DR=0.99
EKF [Sun et al. 2013]	10,000 samples	50% samples, same node	6	0.22	0.42	0.7
Liu et al. [2007]	4,096 nodes	10%–25% nodes	10	0.01	0.01	0.07
ART [Tanachaiwiwat and Helmy 2005]	100 nodes	30%–50% samples, random selection of malicious nodes	100	0.25	0.22	0.21
Chatzigiannakis and Papavassiliou [2007]	40 nodes	10% nodes	40	0.67	0.69	0.7
Chatzigiannakis and Papavassiliou [2007]	40 nodes	40% nodes	40	0.48	0.5	0.6
WTE [Atakli et al. 2008]	100 nodes * 200 samples	0%–25% nodes	100	0.03	0.41	0.78
WTE [Atakli et al. 2008]	400 nodes * 200 samples	0%–25% nodes	400	0.10	0.44	0.78
Bankovic et al. [2010]	2,000 nodes * 2,500 samples (1,000 are used for training)	5% nodes	2,000	0.5	0.5	0.5
Trust-based IDS [Bao et al. 2012]	900 nodes	N/A	N/A	0.001	0.05	N/A
DWE [Oh et al. 2012]	200 samples	20% nodes	20	0.01	0.01	0.02
Dual threshold [Lim and Choi 2013]	100 samples	10% nodes	12	N/A	N/A	0.001
Dual threshold [Lim and Choi 2013]	100 samples	20% nodes	12	0.18	0.14	0.10

equal to each other, apart from noise and errors. The consequence of this assumption is that, by increasing the number of sensors, the information redundancy also increases and the number of sensors taken into account is decisive. Recall from Section 4.4 that the sensing space can be approximately homogeneous only if we consider a small portion of space where there are no obstacles. In works like Chatzigiannakis and Papavassiliou [2007] and Bankovic et al. [2010], where this assumption is not present, the FPR is higher, but the algorithm has wider applicability. Tanachaiwiwat and Helmy [2005] rely on the spatial homogeneity assumption and apply their technique to a large neighborhood (100 nodes). The FPR is better but still not negligible (more than 20%). Atakli et al. [2008] also rely on this assumption and apply their algorithm on very large neighborhoods. With 100 nodes, the FPR for DR = 0.90 is 3%, but for DR = 95 and DR = 99, the FPR increases by an order of magnitude. In contrast, Oh et al. [2012], Bao et al. [2012], and Lim and Choi [2013] are successful in keeping the FPR low even for high DR. Note that with a larger number of nodes the FPR of the technique described

Table V. Detection Performances, Colluding Attacks

Work	Dataset Size	Colluding Percentage	Input Size for Each Algorithm Execution	FPR for DR=0.90	FPR for DR=0.95	FPR for DR=0.99
ART [Tanachaiwiwat and Helmy 2005]	100 nodes	30%–50% samples	100	0.25	0.22	0.21
Chatzigiannakis and Papavassiliou [2007]	40 nodes	10% nodes	40	0.67	0.69	0.7
Chatzigiannakis and Papavassiliou [2007]	40 nodes	40% nodes	40	0.76	0.78	0.8
Robust IF [Rezvani et al. 2013]	20 nodes per 400 samples	40% nodes	20	N/A	0.021	0.021

in Atakli et al. [2008] increases. This result contrasts with the consideration that we made about the the spatial homogeneity assumption. The reason behind that lies probably in the inaccuracy of the empirical ROC curve calculation. Another possible cause is that the algorithm is sensitive to the absolute number of compromised nodes rather than to its ratio to total nodes. For example, 80 out 400 compromised nodes may be harder to detect than 20 out of 100, even though the percentage of malicious nodes is 20% in both cases.

In Table V, we report the results for the cases considering collusion. The results reported in Chatzigiannakis and Papavassiliou [2007] show nonnegligible FPR values (above 60%). The results reported in Tanachaiwiwat and Helmy [2005] have a better FPR (around 20%). Rezvani et al. [2013], instead, achieve very good results (FPR less than 5%). Nevertheless, recall that this technique is applicable only when the spatial homogeneity assumption among the 20 sensors is reasonable. In scenarios where the sensor readings cannot be assumed to share the same physical attribute function, the results may degrade substantially. This is the case for physical attributes like vibration, light, wind, and so forth, where the correlation of the attribute measured at different locations rapidly decreases with the event propagation.

### 7.3. Comparing Techniques' Overhead

The applicability of a technique to a real WSN depends not only on the relationship between the detection rate and the false-positive rate but also on the overhead introduced. We analyze computational and communication overhead for the techniques discussed in the previous section and summarize their asymptotic complexity in Table VI. As usual,  $N$  is the number of sensors,  $N_n$  is the average number of neighbors, and  $W$  is the temporal memory, that is, the number of past samples used.

From Table VI, we note that anomaly detection techniques generally introduce more computational overhead than trust management techniques. The reason behind this result is that trust management iteratively refines its confidence about a sensor's trustworthiness, whereas anomaly detection builds such confidence from scratch at each iteration. On the other hand, this is also the main reason that trust management algorithms are vulnerable to on-off attacks (see Section 5.2).

Another noticeable result is that communication overhead is always kept lower than computational overhead—this result is to be expected since network communication is more expensive in terms of energy and leads to faster battery depletion. In anomaly detection techniques, the communication overhead comes from the

Table VI. Techniques' Overhead

Class	Work	Computational Overhead	Communication Overhead
Anomaly detection	ART [Tanachaiwiwat and Helmy 2005]	$O(W * N_n)$	$O(1)$
	Liu et al. [2007]	$O(N_n^2)$	$O(N_n)$
	Chatzigiannakis and Papavassiliou [2007]	$O(WN_n^2 + N_n^3)$	0
	EKF [Sun et al. 2013]	$O(1)$	$O(N_n)$
	Robust IF [Rezvani et al. 2013]	$O(WN^2)$	0
Trust management	WTE [Atakli et al. 2008]	$O(N_n)$	0
	Bankovic et al. [2010]	$O(N_n^2) + O(W^2)$	0
	Trust-based IDS [Bao et al. 2012]	$O(N_n)$	$O(N_n)$
	DWE [Oh et al. 2012]	$O(N_n)$	0
	Dual threshold [Lim and Choi 2013]	$O(N_n)$	0

execution of consensus-like protocols that decide about the maliciousness of nodes after anomalies are detected. Trust management techniques instead delegate such decisions to the nodes that are higher in a WSN hierarchy (e.g., the forwarding nodes, cluster heads, basestation). Thus, communication overhead is introduced in trust management techniques only when recommendations are enabled (such as in Bao et al. [2012]).

## 8. CONCLUSIONS AND FUTURE RESEARCH DIRECTIONS

Malicious data injections are a considerable threat for WSNs. We reviewed state-of-the-art techniques that can detect malicious data injections by defining an expected behavior and then detecting deviations from it. We classified these approaches into two main families: *anomaly detection* and *trust management*. They differ in the assessment of an anomalous condition, but both rely on the definition of an expected behavior. We analyzed and compared the techniques by their definition of expected behavior and noted that expectations can come from the following correlations: (1) *in time*: different time, same sensor, same attribute; (2) *in space*: same time, different sensors, same attribute; (3) *across different physical attributes*: same time, same sensor, different attributes; or (4) their combination.

While many techniques can be applied, comparatively few target explicitly malicious data injections, especially when collusion between compromised sensors is considered. Most techniques aim to detect erroneous measurements, either to improve the quality of the measuring process (e.g., Subramaniam et al. [2006] and Bettencourt et al. [2007]) or to reduce the power associated with the transmission of the measurements (e.g., Wang et al. [2010] and Salem et al. [2013]).

Work aimed at detecting malicious data injections generally uses spatial correlation in constructing the expectations (e.g., Zhang et al. [2006], Liu et al. [2007], and Chatzigiannakis and Papavassiliou [2007]), in keeping with a general assumption that only a subset of sensors has been compromised. In this case, a nonvoid set of genuine measurements is always present in the spatial domain.

We discussed the different assumptions that characterize the spatial domain and analyzed how they impact the performance of the detection algorithms. More precisely, we observed a substantial decrease in performance when moving away from a homogeneous space model, where all sensors perceive similar measurements, to heterogeneous space models, where different measurements are expected at different locations. This result is visible, for example, in the difference between the results achieved in Tanachaiwiwat and Helmy [2005] and Rezvani et al. [2013], who assume a homogeneous space, and those achieved by Chatzigiannakis and Papavassiliou [2007], who

only assume some degree of correlation between the sensors. The results in the latter case show noticeable higher false-positive rates. We conclude that more research is needed to achieve better results when the spatial domain is heterogeneous. This will also improve the general applicability of the algorithms in real-life deployments.

We explored different approaches to the detection phase, where the deviation from the expected behavior is assessed, and noted a clear preference in the literature for outlier detection techniques (e.g., Ngai et al. [2006], Liu et al. [2007], and Sun et al. [2013]). In this case, the expectation of a measurement is compliant with a generalization of the measurement's behavior. This approach is independent from the context and is preferred to more context-specific techniques based on model checking (e.g., Handschin et al. [1975]).

Finally, to complete the detection of malicious data injections, we identified two main aspects that need to be addressed: *diagnosis* and *characterization*. These are, by and large, insufficiently studied in the literature.

Diagnosis consists of identifying the cause of the detected anomaly, which, besides malicious data injections, may lie in faults or events of interest. Both these phenomena can produce deviations from expected behavior similar to malicious injections. While partial diagnosis is investigated in, for example, Tanachaiwiwat and Helmy [2005], Bettencourt et al. [2007], Chatzigiannakis and Papavassiliou [2007], and Oh et al. [2012], an exhaustive diagnosis phase is still lacking. Fault-related anomalies may be handled separately from malicious data injections, as fault models are relatively well categorized and understood. However, event-related anomalies cannot be considered separately (as in Liu et al. [2007]), since an attacker may inject malicious measurements that depict a fabricated event or conceal a real event. Therefore, in WSNs that monitor the occurrence of events, malicious injections and events should be addressed together, to produce a compromise-resistant detection and characterization of events.

Similarly, further investigation of the *attack characterization* is needed, in particular to identify the compromised sensors in the presence of collusion. This aspect adds more complexity to the problem since colluding sensors can reduce data inconsistencies introduced in the attack, especially in the spatial domain.

Across all of the aspects, a good model of expected system behavior plays a central role and determines both the applicability of the algorithms for detecting malicious data injections and their performance.

## ACKNOWLEDGMENTS

We wish to thank the anonymous reviewers for their comments and suggestions, which have made a valuable contribution to this paper. We are grateful to Dr. Igor Muttik from Intel and to the colleagues in our research group, who have contributed to this work through many useful discussions. This work was funded as part of the Intel Collaborative Research Institute on Sustainable Connected Cities.

## REFERENCES

- Idris M. Atakli, Hongbing Hu, Yu Chen, Wei-Shinn Ku, and Zhou Su. 2008. Malicious node detection in wireless sensor networks using weighted trust evaluation. In *SpringSim*, Hassan Rajaei, Gabriel A. Wainer, and Michael J. Chinni (Eds.). SCS/ACM, 836–843.
- Majid Bahrepour, Yang Zhang, Nirvana Meratnia, and Paul J. M. Havinga. 2009. Use of event detection approaches for outlier detection in wireless sensor networks. In *2009 5th International Conference on Intelligent Sensors, Sensor Networks and Information Processing (ISSNIP'09)*. IEEE, 439–444.
- Zorana Bankovic, José Manuel Moya, Álvaro Araujo, David Fraga, Juan Carlos Vallejo, and Juan-Mariano de Goyeneche. 2010. Distributed intrusion detection system for wireless sensor networks based on a reputation system coupled with kernel self-organizing maps. *Integrated Computer-Aided Engineering* 17, 2, 87–102.

- Fenye Bao, Ing-Ray Chen, Moonjeong Chang, and Jin-Hee Cho. 2012. Hierarchical trust management for wireless sensor networks and its applications to trust-based routing and intrusion detection. *IEEE Transactions on Network and Service Management* 9, 2, 169–183.
- Elisa Bertino, Aleksandar Ignatovic, and Sanjay Jha. 2014. Secure data aggregation technique for wireless sensor networks in the presence of collusion attacks. *IEEE Transactions on Dependable and Secure Computing* 99, PrePrints, 1. DOI : <http://dx.doi.org/10.1109/TDSC.2014.2316816>
- Luís M. A. Bettencourt, Aric A. Hagberg, and Levi B. Larkey. 2007. Separating the wheat from the chaff: Practical anomaly detection schemes in ecological applications of distributed sensor networks. In *DCOSS (Lecture Notes in Computer Science)*, James Aspnes, Christian Scheideler, Anish Arora, and Samuel Madden (Eds.), Vol. 4549. Springer, 223–239.
- Azzedine Boukerche. 2009. *Algorithms and Protocols for Wireless Sensor Networks*. Wiley-IEEE Press.
- Azzedine Boukerche, Horacio A. B. F Oliveira, Eduardo F. Nakamura, and Antonio A. F. Loureiro. 2008. Secure localization algorithms for wireless sensor networks. *IEEE Communications Magazine* 46, 4 (2008), 96–101.
- Levente Buttyan and Jean-Pierre Hubaux. 2008. *Security and Cooperation in Wireless Networks: Thwarting Malicious and Selfish Behavior in the Age of Ubiquitous Computing*. Cambridge University Press.
- Claude Castelluccia, Aurélien Francillon, Daniele Perito, and Claudio Soriente. 2009. On the difficulty of software-based attestation of embedded devices. In *ACM Conference on Computer and Communications Security*, Ehab Al-Shaer, Somesh Jha, and Angelos D. Keromytis (Eds.). ACM, 400–409.
- Varun Chandola, Arindam Banerjee, and Vipin Kumar. 2009. Anomaly detection: A survey. *ACM Computer Surveys* 41, 3.
- Vassilis Chatzigiannakis and Symeon Papavassiliou. 2007. Diagnosing anomalies and identifying faulty nodes in sensor networks. *IEEE Sensors Journal* 7, 5 (2007), 637–645.
- Ethan W. Dereszynski and Thomas G. Dietterich. 2011. Spatiotemporal models for data-anomaly detection in dynamic environmental monitoring campaigns. *TOSN* 8, 1, 3.
- Wenliang Du, Jing Deng, Yunghsiang S. Han, Pramod K. Varshney, Jonathan Katz, and Aram Khalili. 2005. A pairwise key predistribution scheme for wireless sensor networks. *ACM Transactions on Information Systems Security* 8, 2, 228–258.
- Elena Fasolo, Michele Rossi, Jörg Widmer, and Michele Zorzi. 2007. In-network aggregation techniques for wireless sensor networks: A survey. *IEEE Wireless Communications* 14, 2, 70–87.
- Saurabh Ganeriwal, Laura Balzano, and Mani B. Srivastava. 2003. Reputation-based framework for high integrity sensor networks. *TOSN* 4, 3.
- Saurabh Ganeriwal and Mani B. Srivastava. 2004. Reputation-based framework for high integrity sensor networks. In *SASN*, Sanjeev Setia and Vipin Swarup (Eds.). ACM, 66–77.
- Shuo Guo, Ziguo Zhong, and Tian He. 2009. FIND: Faulty node detection for wireless sensor networks. In *SenSys*, David E. Culler, Jie Liu, and Matt Welsh (Eds.). ACM, 253–266.
- Edmund Handschin, Fred C. Schweppe, Jurg Kohlas, and Armin Fiechter. 1975. Bad data analysis for power system state estimation. *IEEE Transactions on Power Apparatus and Systems* 94, 2 (1975), 329–337.
- Wendi Rabiner Heinzelman, Anantha Chandrakasan, and Hari Balakrishnan. 2000. Energy-efficient communication protocol for wireless microsensor networks. In *HICSS*.
- Cheryl V. Hinds. 2009. Efficient detection of compromised nodes in a wireless sensor network. In *SpringSim*, Gabriel A. Wainer, Clifford A. Shaffer, Robert M. McGraw, and Michael J. Chinni (Eds.). SCS/ACM.
- Lei Huang, Lei Li, and Qiang Tan. 2006. Behavior-based trust in wireless sensor network. In *APWeb Workshops (Lecture Notes in Computer Science)*, Heng Tao Shen, Jinbao Li, Minglu Li, Jun Ni, and Wei Wang (Eds.), Vol. 3842. Springer, 214–223.
- J. Edward Jackson and Govind S. Mudholkar. 1979. Control procedures for residuals associated with principal component analysis. *Technometrics* 21, 3, 341–349.
- Dharanipragada Janakiram, Vanteddu Adi Mallikarjuna Reddy, and A. V. U. Phani Kumar. 2006. Outlier detection in wireless sensor networks using Bayesian belief networks. In *First International Conference on Communication System Software and Middleware (Comsware'06)*. 1–6.
- Raja Jurdak, X. Rosalind Wang, Oliver Obst, and Philip Valencia. 2011. Wireless sensor network anomalies: Diagnosis and detection strategies. In *Intelligence-Based Systems Engineering*, Andreas Tolk and Lakhmi C. Jain (Eds.). Intelligent Systems Reference Library, Vol. 10. Springer, Berlin, 309–325. DOI : [http://dx.doi.org/10.1007/978-3-642-17931-0\\_12](http://dx.doi.org/10.1007/978-3-642-17931-0_12)
- Rudolph Emil Kalman. 1960. A new approach to linear filtering and prediction problems. *Transactions of the ASME-Journal of Basic Engineering* 82, Series D, 35–45.
- Chris Karlof and David Wagner. 2003. Secure routing in wireless sensor networks: Attacks and countermeasures. *Ad Hoc Networks* 1, 2–3, 293–315.

- Muhammad Khurram Khan and Khaled Alghathbar. 2010. Cryptanalysis and security improvements of 'two-factor user authentication in wireless sensor networks'. *Sensors* 10, 3, 2450–2459. DOI : <http://dx.doi.org/10.3390/s100302450>
- Pavel Laskov, Christin Schäfer, Igor V. Kotenko, and Klaus-Robert Müller. 2004. Intrusion detection in unlabeled data with quarter-sphere support vector machines. *Praxis der Informationsverarbeitung und Kommunikation* 27, 4, 228–236.
- Sung Yul Lim and Yoon-Hwa Choi. 2013. Malicious node detection using a dual threshold in wireless sensor networks. *Journal of Sensor and Actuator Networks* 2, 1, 70–84.
- An Liu and Peng Ning. 2008. TinyECC: A configurable library for elliptic curve cryptography in wireless sensor networks. In *IPSN*. IEEE Computer Society, 245–256.
- Fang Liu, Xiuzhen Cheng, and Dechang Chen. 2007. Insider attacker detection in wireless sensor networks. In *INFOCOM*. IEEE, 1937–1945.
- Javier Lopez, Rodrigo Roman, Isaac Agudo, and M. Carmen Fernández Gago. 2010. Trust management systems for wireless sensor networks: Best practices. *Computer Communications* 33, 9, 1086–1093.
- Xuanwen Luo, Ming Dong, and Yinlun Huang. 2006. On distributed fault-tolerant detection in wireless sensor networks. *IEEE Transactions on Computers* 55, 1, 58–70.
- Stephen Marsland. 2009. *Machine Learning - An Algorithmic Perspective*. CRC Press. I–XVI, 1–390 pages.
- John McHugh. 2000. Testing intrusion detection systems: A critique of the 1998 and 1999 DARPA intrusion detection system evaluations as performed by Lincoln laboratory. *ACM Transactions on Information System Security* 3, 4, 262–294.
- Mohammad Momani, Subhash Challa, and Rami Alhmouz. 2008. Can we trust trusted nodes in wireless sensor networks? In *International Conference on Computer and Communication Engineering (ICCCE'08)*. IEEE, 1227–1232.
- Edith C. H. Ngai, Jiangchuan Liu, and Michael R. Lyu. 2006. On the intruder detection for sinkhole attack in wireless sensor networks. In *IEEE International Conference on Communications (ICC'06)*. Vol. 8. IEEE, 3383–3389.
- Kevin Ni and Gregory J. Pottie. 2012. Sensor network data fault detection with maximum a posteriori selection and bayesian modeling. *TOSN* 8, 3, 23.
- Seo Hyun Oh, Chan O. Hong, and Yoon-Hwa Choi. 2012. A malicious and malfunctioning node detection scheme for wireless sensor networks. *Wireless Sensor Network* 4, 3, 84–90.
- Suat Özdemir and Yang Xiao. 2009. Secure data aggregation in wireless sensor networks: A comprehensive overview. *Computer Networks* 53, 12, 2022–2037.
- Spiros Papadimitriou, Hiroyuki Kitagawa, Phillip B. Gibbons, and Christos Faloutsos. 2003. LOCI: Fast outlier detection using the local correlation integral. In *ICDE*, Umeshwar Dayal, Krithi Ramamritham, and T. M. Vijayaraman (Eds.). IEEE Computer Society, 315–326.
- Lilia Paradis and Qi Han. 2007. A survey of fault management in wireless sensor networks. *Journal of Network System Management* 15, 2, 171–190.
- Taejoon Park and Kang G. Shin. 2005. Soft tamper-proofing via program integrity verification in wireless sensor networks. *IEEE Transactions on Mobile Computing* 4, 3, 297–309.
- Adrian Perrig, John Stankovic, and David Wagner. 2004. Security in wireless sensor networks. *Communications of the ACM*, 53–57. Issue 6. DOI : <http://dx.doi.org/10.1145/990680.990707>
- Bartosz Przydatek, Dawn Xiaodong Song, and Adrian Perrig. 2003. SIA: Secure information aggregation in sensor networks. In *SenSys*, Ian F. Akyildiz, Deborah Estrin, David E. Culler, and Mani B. Srivastava (Eds.). ACM, 255–265.
- Sutharshan Rajasegarar, James C. Bezdek, Christopher Leckie, and Marimuthu Palaniswami. 2009. Elliptical anomalies in wireless sensor networks. *TOSN* 6, 1.
- Sutharshan Rajasegarar, Christopher Leckie, and Marimuthu Palaniswami. 2008. Anomaly detection in wireless sensor networks. *IEEE Wireless Communications* 15, 4, 34–40.
- Sutharshan Rajasegarar, Christopher Leckie, Marimuthu Palaniswami, and James C. Bezdek. 2006. Distributed anomaly detection in wireless sensor networks. In *10th IEEE Singapore International Conference on Communication Systems (ICCS'06)*. IEEE, 1–5.
- Sutharshan Rajasegarar, Christopher Leckie, Marimuthu Palaniswami, and James C. Bezdek. 2007. Quarter sphere based distributed anomaly detection in wireless sensor networks. In *ICC (2009-04-15)*. IEEE, 3864–3869.
- Murad A. Rassam, Anazida Zainal, and Mohd Aizaini Maarof. 2013. Advancements of data anomaly detection research in wireless sensor networks: A survey and open issues. *Sensors* 13, 8, 10087–10122. DOI : <http://dx.doi.org/10.3390/s130810087>

- Maxim Raya, Panos Papadimitratos, Virgil D. Gligor, and Jean-Pierre Hubaux. 2008. On data-centric trust establishment in ephemeral ad hoc networks. In *27th IEEE Conference on Computer Communications (INFOCOM'08)*.
- Mohsen Rezvani, Aleksandar Ignjatovic, Elisa Bertino, and Sanjay Jha. 2013. A robust iterative filtering technique for wireless sensor networks in the presence of malicious attacks. In *SenSys*, Chiara Petrioli, Landon P. Cox, and Kamin Whitehouse (Eds.). ACM, 30.
- John A. Rice. 2007. *Mathematical Statistics and Data Analysis* (3rd ed.). Duxbury Press.
- Sandip Roy, Marco Conti, Sanjeev Setia, and Sushil Jajodia. 2014. Secure data aggregation in wireless sensor networks: Filtering out the attacker's impact. *IEEE Transactions on Information Forensics and Security* 9, 4 (2014), 681–694.
- Osman Salem, Yaning Liu, and Ahmed Mehaoua. 2013. A lightweight anomaly detection framework for medical wireless sensor networks. In *WCNC*. IEEE, 4358–4363.
- Yingpeng Sang, Hong Shen, Yasushi Inoguchi, Yasuo Tan, and Naixue Xiong. 2006. Secure data aggregation in wireless sensor networks: A survey. In *PDCAT*. IEEE Computer Society, 315–320.
- Arvind Seshadri, Mark Luk, Adrian Perrig, Leendert van Doorn, and Pradeep K. Khosla. 2006. SCUBA: Secure code update by attestation in sensor networks. In *Workshop on Wireless Security*, Radha Poovendran and Ari Juels (Eds.). ACM, 85–94.
- Arvind Seshadri, Adrian Perrig, Leendert van Doorn, and Pradeep K. Khosla. 2004. SWATT: Software-based attestation for embedded devices. In *IEEE Symposium on Security and Privacy*. IEEE Computer Society, 272.
- Nauman Shahid, Ijaz Haider Naqvi, and Saad B. Qaisar. 2012. Quarter-sphere SVM: Attribute and spatio-temporal correlations based outlier & event detection in wireless sensor networks. In *WCNC*. IEEE, 2048–2053.
- Abhishek B. Sharma, Leana Golubchik, and Ramesh Govindan. 2010. Sensor faults: Detection methods and prevalence in real-world datasets. *TOSN* 6, 3.
- Shigen Shen, Guangxue Yue, Qiyang Cao, and Fei Yu. 2011. A survey of game theory in wireless sensor networks security. *JNW* 6, 3, 521–532.
- Timothy J. Shepard. 1996. A channel access scheme for large dense packet radio networks. In *SIGCOMM*. 219–230. <http://dblp.uni-trier.de/db/conf/sigcomm/sigcomm1996.html#Shepard96>.
- Suresh Singh, Mike Woo, and C. S. Raghavendra. 1998. Power-aware routing in mobile ad hoc networks. In *MOBICOM*, William P. Osborne and Dhawal B. Moghe (Eds.). ACM, 181–190.
- Sharmila Subramaniam, Themis Palpanas, Dimitris Papadopoulos, Vana Kalogeraki, and Dimitrios Gunopulos. 2006. Online outlier detection in sensor data using non-parametric models. In *VLDB*, Umeshwar Dayal, Kyu-Young Whang, David B. Lomet, Gustavo Alonso, Guy M. Lohman, Martin L. Kersten, Sang Kyun Cha, and Young-Kuk Kim (Eds.). ACM, 187–198.
- Bo Sun, Xuemei Shan, Kui Wu, and Yang Xiao. 2013. Anomaly detection based secure in-network aggregation for wireless sensor networks. *IEEE Systems Journal* 7, 1, 13–25.
- Yan Sun, Hong Luo, and Sajal K. Das. 2012. A trust-based framework for fault-tolerant data aggregation in wireless multimedia sensor networks. *IEEE Transactions on Dependable and Secure Computing* 9, 6, 785–797.
- Yan Lindsay Sun, Zhu Han, Wei Yu, and K. J. Ray Liu. 2006. A trust evaluation framework in distributed networks: Vulnerability analysis and defense against attacks. In *INFOCOM*. IEEE.
- Sapon Tanachaiwiwat and Ahmed Helmy. 2005. Correlation analysis for alleviating effects of inserted data in wireless sensor networks. In *MobiQuitous*. IEEE Computer Society, 97–108.
- Xue Wang, Liang Ding, and Daowei Bi. 2010. Reputation-enabled self-modification for target sensing in wireless sensor networks. *IEEE Transactions on Instrumentation and Measurement* 59, 1, 171–179.
- Weili Wu, Xiuzhen Cheng, Min Ding 0001, Kai Xing, Fang Liu, and Ping Deng. 2007. Localized outlying and boundary data detection in sensor networks. *IEEE Transactions on Knowledge and Data Engineering* 19, 8, 1145–1157.
- Miao Xie, Song Han, Biming Tian, and Sazia Parvin. 2011. Anomaly detection in wireless sensor networks: A survey. *Journal on Network and Computer Applications* 34, 4, 1302–1325.
- Yi Yang, Xinran Wang, Sencun Zhu, and Guohong Cao. 2006. SDAP: A secure hop-by-hop data aggregation protocol for sensor networks. In *MobiHoc*, Sergio Palazzo, Marco Conti, and Raghupathy Sivakumar (Eds.). ACM, 356–367.
- Yanli Yu, Keqiu Li, Wanlei Zhou, and Ping Li. 2012. Trust mechanisms in wireless sensor networks: Attack analysis and countermeasures. *Journal on Network and Computer Applications* 35, 3, 867–880. DOI: <http://dx.doi.org/10.1016/j.jnca.2011.03.005>



- Theodore Zahariadis, Helen-Catherine Leligou, Panagiotis Trakadas, and Stamatis Voliotis. 2010a. Mobile networks trust management in wireless sensor networks. *European Transactions on Telecommunications* 21, 4, 386–395.
- Theodore Zahariadis, Helen C. Leligou, Panagiotis Trakadas, and Stamatis Voliotis. 2010b. Trust management in wireless sensor networks. *European Transactions on Telecommunications* 21, 4, 386–395. DOI: <http://dx.doi.org/10.1002/ett.1413>
- Dazhi Zhang and Donggang Liu. 2010. DataGuard: Dynamic data attestation in wireless sensor networks. In *DSN*. IEEE, 261–270.
- Wei Zhang, Sajal K. Das, and Yonghe Liu. 2006. A trust based framework for secure data aggregation in wireless sensor networks. In *SECON*. IEEE, 60–69.
- Yang Zhang, Nicholas A. S. Hamm, Nirvana Meratnia, Alfred Stein, M. van de Voort, and Paul J. M. Havinga. 2012. Statistics-based outlier detection for wireless sensor networks. *International Journal of Geographical Information Science* 26, 8 (2012), 1373–1392.

Received November 2014; revised May 2015; accepted July 2015