

A Framework for Identifying Compromised Nodes in Wireless Sensor Networks

QING ZHANG, TING YU and PENG NING
North Carolina State University

Sensor networks are often subject to physical attacks. Once a node's cryptographic key is compromised, an attacker may completely impersonate it, and introduce arbitrary false information into the network. Basic cryptographic mechanisms are often not effective in this situation. Most techniques to address this problem focus on detecting and tolerating false information introduced by compromised nodes. They cannot pinpoint exactly where the false information is introduced and who is responsible for it.

In this paper, we propose an application-independent framework for accurately identifying compromised sensor nodes. The framework provides an appropriate abstraction of application-specific detection mechanisms, and models the unique properties of sensor networks. Based on the framework, we develop alert reasoning algorithms to identify compromised nodes. The algorithm assumes that compromised nodes may collude at will. We show that our algorithm is optimal in the sense that it identifies the largest number of compromised nodes without introducing false positives. We evaluate the effectiveness of the designed algorithm through comprehensive experiments.

Categories and Subject Descriptors: C.2.0 [Computer-Communication Networks]: General—*Security and protection*; K.6.5 [Management of Computing and Information Systems]: Security and Protection

General Terms: Algorithm, Security

Additional Key Words and Phrases: Sensor network, intrusion detection

1. INTRODUCTION

Compared with traditional wired and wireless networks, low-power wireless sensor networks can be rapidly deployed in a large geographical area in a self-configured manner. This makes them particularly suitable for real-time, large-scale information collection and event monitoring for mission-critical applications in hostile environments, such as target tracking and battlefield surveillance.

The work of Zhang and Yu was partially supported by the National Science Foundation (NSF) under grants IIS-0430274. Ning's work was supported by the NSF under grants CAREER-0447761 and CNS-0430223.

A preliminary version of this paper appeared in *Proceedings of 2nd IEEE Communications Society/CreateNet International Conference on Security and Privacy in Communication Networks (SecureComm 2006)*, August 2006 [Zhang et al. 2006].

Authors' address: Cyber Defense Laboratory, Department of Computer Science, North Carolina State University, Raleigh, NC 27695; emails: {qzhang4.tyu.pning}@ncsu.edu

Permission to make digital/hard copy of all or part of this material without fee for personal or classroom use provided that the copies are not made or distributed for profit or commercial advantage, the ACM copyright/server notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists requires prior specific permission and/or a fee.

© 20YY ACM 0000-0000/20YY/0000-0001 \$5.00

Such applications, meanwhile, impose unique security challenges. Since sensors are extremely resource-constrained, many existing security mechanisms cannot be directly applied to sensor networks [Perrig et al. 2001]. In recent years, we have witnessed great efforts toward designing security primitives specific for wireless sensor networks. For example, a variety of authentication and key management schemes have been proposed and implemented [Camtepe and Yener 2004; Liu et al. 2004]. On the other hand, since sensors are often deployed in open environments, they are vulnerable to physical attacks. Once recovering the keying materials of some nodes, an adversary is able to impersonate them completely, and inject arbitrary false information. Basic cryptographic mechanisms, such as authentication and integrity protection, are usually not effective against such impersonation attacks [Du et al. 2005].

Recently, several approaches have been proposed to cope with compromised nodes. These approaches mainly fall into two categories. Approaches in the first category are to detect and tolerate false information introduced by attackers [Du et al. 2003a; Przydatek et al. 2003; Hu and Evans 2003; Ye et al. 2004], in particular during data aggregation. Once the base station receives aggregated data, it checks their validity through mechanisms such as sampling and deployment of redundant sensors. However, these techniques often cannot be used to identify where the false information is introduced and who is responsible for it.

Approaches in the second category rely on application specific detection mechanisms which enable sensor nodes to monitor the activities of others nearby. Once an abnormal activity is observed, a node may raise an alert either to the base station or to other nodes, who further determine which nodes are compromised. We call approaches in this category *alert-based*. Representative alert-based approaches include those in sensor network routing [Ganerwal and Srivastava 2004] and localization [Liu et al. 2005].

Alerts from sensor nodes make it possible to pinpoint compromised nodes. However, how to effectively utilize such information is a very challenging problem. It is hard to decide whether an alert can be trusted since it is very likely that compromised nodes raise false alerts to mislead the base station and other nodes. Compromised nodes may further form a local majority in the network, and collude, increasing their influences in the network. Further, existing alert-based approaches are specific to certain applications, and cannot be easily extended to other domains. A general solution to the *accurate identification* of compromised nodes still remain elusive.

The problem of identifying compromised nodes shares certain similarity with fault diagnosis in diagnosable systems [Preparata et al. 1967; Araki and Shibata 2003; Sullivan 1988; Dahbura and Masson 1984; Fuhrman 1996]. However, in those systems, faults are assumed to be permanent, which means a faulty node will always fail a test, and thus can always be identified by fault-free nodes. Some later works relax permanent faults to intermittent faults [Dahbura et al. 1987; Kozłowski and Krawczyk 1991], which however still assume that a faulty node cannot pass a test following certain probabilities. These assumptions do not hold in sensor networks, where a compromised node may behave arbitrarily. For example, it may always report correct sensing data, and meanwhile issue false alerts. Such malicious behavior cannot be observed by an uncompromised node. Thus, we cannot directly apply works in self-diagnosable systems to identify compromised nodes in sensor networks.

The problem of false alerts (or feedback) and collusion from malicious entities also arises in other decentralized systems such as online auction communities and P2P sys-

tems [Aberer and Despotovic 2001; Kamvar et al. 2003; Lee et al. 2003; Mui et al. 2002; Richardson et al. 2003; Xiong and Liu 2002; Yu and Singh 2002]. Reputation-based trust management has been adopted as an effective means to form cooperative groups in the above systems. One seemingly attractive approach is to apply existing trust management techniques in sensor networks. For example, we may identify sensor nodes with the lowest trust values as compromised. However, as a decentralized system, sensor networks bear quite unique properties, and significantly differ from the above systems. Many assumptions in reputation-based trust management do not hold in sensor networks. Thus, simply applying those techniques is unlikely to be effective (see section 4 for a detailed experimental comparison).

For example, in P2P systems, interactions may happen between any two entities. If an entity provides misleading information or poor services, it is likely that some other entities will be able to detect it and issue negative feedback accordingly. The interactions between sensor nodes, however, are restricted by the deployment of a sensor network. For a given node, only a fixed set of nodes are able to observe it. Thus, it is easy for compromised nodes to form local majorities.

Also, most decentralized environments are composed of autonomous entities, which pursue to maximize their own interests. Incentive mechanisms are needed to encourage entities to issue (or at least not discourage them from issuing) feedback about others. A sensor network, on the other hand, is essentially a distributed system, where all the sensor nodes are designed to cooperate and finish a common task. Therefore, it is possible to design identification mechanisms that achieve global optimality for a given goal. For instance, to cope with false alerts, we may choose to identify as compromised both the target and the issuer of an alert, as long as it will improve the security of the whole system. Such an approach is usually not acceptable in P2P systems and online auction communities.

Indeed, the unique properties of sensor networks bring both challenges and opportunities. How to accommodate and take advantages of these properties is the key to the accurate identification of compromised nodes.

In this paper, we propose novel techniques to provide general solutions to the identification of compromised sensor nodes. Our techniques are based on an application-independent framework that abstracts some of the unique and intrinsic properties of sensor networks. Therefore, it can be used to model a large range of existing sensor network applications. In summary, the contributions of this paper include the following:

- (1) We develop an application-independent framework for identifying compromised nodes based on alerts generated by specific detection mechanisms in sensor networks. The central component of the framework is an abstraction of the monitoring relationship between sensor nodes. Such relationship can be derived from application specific detection mechanisms. The framework further models sensor nodes' sensing and monitoring capabilities and their impacts on detection accuracy. We show by example that many existing sensor networks can be easily modeled by the proposed framework. As our framework is built on the alert-based detection mechanisms provided by applications, it does not require sensor nodes to support additional functionalities, nor does it impose additional communication and computation costs to the network.
- (2) Based on the proposed framework, we design an alert reasoning algorithm to accurately identify compromised sensor nodes. The algorithm does not rely on any assumptions on how compromised nodes behave and collude. We show that the algo-

rithm is optimal, in the sense that given any set of alerts, our algorithm identifies the largest number of compromised nodes that can generate these alerts, without introducing any false positives. We also study how to tradeoff certain false positives to further eliminate compromised nodes.

- (3) To better understand the capability of the above reasoning algorithm, we further consider a special case where all the compromised nodes actually collude with each other (i.e., they behave consistently and do not raise alerts against each other). We develop an identification algorithm when assuming the base station is aware of such collusion. The identification capability of this special algorithm serves as an upper bound for that of the general algorithm.
- (4) We conduct comprehensive experiments to evaluate the proposed algorithm. The results show that it yields high detection rates and bounded false positive rates, and thus is effective in identifying compromised nodes.

The rest of the paper is organized as follows. Section 2.3 presents a general framework for identifying compromised nodes, and shows how sensor network application can be modeled by the framework. In section 3, we present algorithms that identify compromised sensor nodes with optimal accuracy, for both the collusion and non-collusion cases. In section 4, we show the effectiveness of our algorithms through experimental evaluation. Some relevant issues are discussed in section 5. Section 6 reports closely related work to this paper. Concluding remarks are given in section 7.

2. A GENERAL FRAMEWORK FOR IDENTIFYING COMPROMISED NODES

In this section, we first use an example to identify the aspects of sensor networks that are relevant to the identification of compromised nodes. We then present the general framework.

2.1 An Example Sensor Network Application

Many sensor network applications require sensors' location information (e.g., in target tracking). Since it is often too expensive to equip localization devices such as GPS receivers on every node, many location discovery algorithms depend on beacon nodes, i.e., sensor nodes aware of their locations. A non-beacon node queries beacon nodes nearby for references, and estimates its own location. An example deployment of the sensor network is shown in figure 1, where beacon nodes and non-beacon nodes are represented by square and round nodes respectively. An edge from a beacon node b to a non-beacon node s indicates that b can provide location references to s .

A compromised beacon node may claim its own location arbitrarily, making non-beacon nodes around derive their locations incorrectly. Liu et al. [Liu et al. 2005] proposed a mechanism to detect malicious beacon nodes. The basic idea is to let beacon nodes probe each other and check the sanity of the claimed locations. Suppose beacon node b_1 's location is (x, y) and beacon node b_2 claims its location to be (x', y') . If the difference between the derived distance and the measured distance exceeds a threshold ϵ , then b_1 will consider b_2 as compromised, and report to the base station. Clearly, a compromised beacon node may also send false alerts to the base station.

After receiving a set of alerts from beacon nodes, what information is needed by the base station to make a rational decision on compromised nodes? First, the base station has to know whether an alert is valid, i.e., whether beacon nodes b_1 and b_2 are close enough so

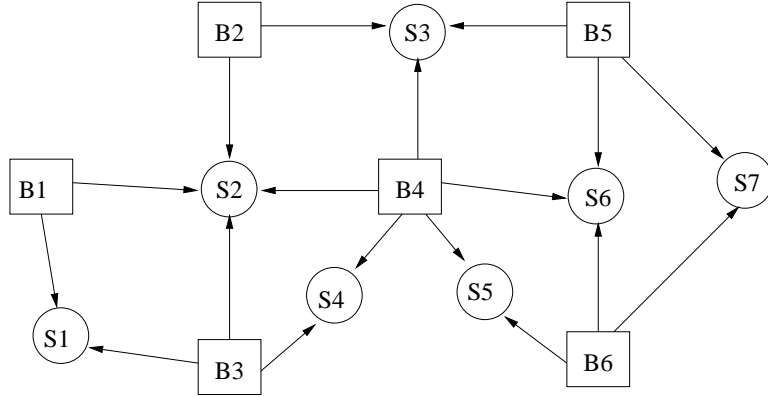


Fig. 1. The deployment of beacon nodes in sensor network localization

that they can probe each other. In other words, the monitoring relationship between beacon nodes is needed.

Second, due to the imprecision of distance measuring, it is possible that an uncompromised beacon node raises an alert against another uncompromised one. The base station has to take this possibility into consideration.

Third, it is necessary to regularly probe a beacon node so that it can be detected promptly if that beacon node is compromised and provides misleading location references.

The above information is quite relevant for the base station to reason about compromised nodes, and is commonly available in sensor network applications. Thus, it should be included by a general framework for identifying compromised nodes.

2.2 Assumptions

We make the following assumptions about sensor network applications before presenting a general framework for identifying compromised nodes.

First, we assume there exist application-specific detection mechanisms deployed in a sensor network, which enable sensor nodes to observe each other's behavior. Such detection mechanisms are commonly employed in sensor networks. Examples include beacon node probing in sensor network localization as mentioned above, witnesses in data aggregation [Du et al. 2003a] and the watchdog mechanism [Ganerwal and Srivastava 2004]. A sensor node s_1 is called an *observer* of another node s_2 if s_1 can observe the behavior of s_2 . A node may have multiple observers or observe multiple other nodes.

The detection mechanisms are not assumed to be completely accurate. But we do assume that, given a sufficient number of observations of the same node, the number of alerts issued by two uncompromised observers should be statistically consistent. In other words, if a node's behavior during a time period T is abnormal, then all the uncompromised observers should raise a significantly high number of alerts during T . Whether the number of alerts is significantly high depends on the characteristics of the sensor node and the detection mechanism, which will be discussed later in section 2.3.

Second, we focus on static sensor networks, where sensor nodes do not change their locations dramatically once deployed. A large range of sensor networks fall into this cate-

gory, e.g., target tracking and environmental monitoring. One consequence of this assumption is that the observability relationship between sensor nodes does not change unless a sensor network is reconfigured.

Third, we assume that message confidentiality and integrity are protected through key management and authentication mechanisms [Du et al. 2003b], so that a sensor node can send information securely and reliably to the base station. Several techniques have been proposed in the literature to ensure the availability of such channels [Bose et al. 2001; Deng et al. 2004].

Finally, we assume the base station of a sensor network is trusted, and has sufficient computation and communication capabilities. Hence, we adopt a centralized approach, where the base station is responsible for reasoning about the alerts and identifying compromised nodes. The responsibility of each node is only to observe abnormal activities and raise alerts to the base station. We will briefly discuss in section 5 decentralized approaches, where sensor nodes also take part in the reasoning and identification process.

2.3 The Framework

With the above assumptions, a general framework for identifying compromised nodes is composed of the following components:

Observability graph. An observability graph is a directed graph $G(V, E)$, where V is a set of vertices that represent sensor nodes, and E is a set of edges. An edge $(s_1, s_2) \in E$ if and only if s_1 is an observer of s_2 . An observability graph is derived from the detection mechanism of an application. V only contains those nodes whose security is concerned, and is involved in the underlying detection mechanism. For example, in the sensor network localization problem, the observability graph only includes beacon nodes.

Alerts. An alert takes the form (t, s_1, s_2) , indicating that node s_1 observes an abnormal activity of s_2 at time t . The information in an alert may be further enriched, for example, by including s_1 's confidence on the alert. For simplicity, we omit such parameters in this paper. Note that alerts may not need to be explicitly sent by sensor nodes. Instead, in some applications they can be implicitly inferred by the base station from the sensing data sent by sensor nodes (e.g., when the reported data of two adjacent sensor nodes differ than a certain threshold).

Sensor behavior model. Sensors are not perfect. Even if a node is uncompromised, it may still occasionally report inaccurate information or behave abnormally. A sensor behavior model includes a parameter r_m that represents the percentage of normal activities conducted by an uncompromised node. We call r_m the *reliability* of sensors. For example, in sensor network localization, if $r_m = 0.99$, then 99% of the time, an uncompromised beacon node provides the correct location references.

Observer model. Similarly, an observer model represents the effectiveness of the detection mechanism of a sensor network, which is captured by its observability rate $r_b(i, j)$, positive accuracy r_p and negative accuracy r_n . Suppose s_1 is an observer of s_2 . $r_b(1, 2)$ is the probability that s_1 observes an activity conducted by s_2 . This reflects the fact that in some applications, due to cost, energy concerns, and distance between nodes, s_1 may not be able to observe every activity of s_2 . For simplicity, in this paper we use a global r_b value between all pair of nodes. The positive accuracy r_p is the probability that s_1 raises an alert when s_2 conducts an abnormal activity observed by s_1 . Similarly,

r_n is the probability that s_1 does not raise an alert when s_2 conducts a normal activity observed by s_1 . r_p and r_n reflect the intrinsic capability of a detection mechanism.

The sensor behavior model and the observer model can usually be obtained from the specification of sensors and an application's detection mechanisms.

Security estimation. If it is possible that all the nodes in the network are compromised, then the base station cannot identify definitely which nodes are compromised based on alerts. Therefore, this framework focuses on the situation where the number of compromised nodes does not exceed a certain threshold K . We call K the *security estimation* of a network. How to determine K is application specific, depending on, e.g., the estimation of attackers' capability, the strength of sensors' keys, and how long the network has been deployed. We emphasize that K is only an *upper bound* of the number of compromised nodes. The base station does not need to know the exact number of compromised nodes in a network.

Identification function. An identification function F determines which nodes are compromised. Formally, it takes as inputs the observability graph G , the sensor reliability r_m , the observer model (r_b, r_p, r_n) , the security estimation K , and a set of alerts raised during a period T , and returns a set of node IDs, which indicate those nodes that are considered compromised.

We note that our framework is built on the alert-based detection mechanisms provided by applications. The framework itself does not require sensor nodes to support additional functionalities, and thus does not introduce additional communication and computation costs to the network.

Further, for simplicity, the above framework assumes that all the nodes and observers follow the same sensor behavior model and observe model. In some applications, there may be different types of sensor nodes with different sensing and observing capabilities. Our framework can be easily extended to model such applications by specifying different models for different types of sensor nodes.

With the above framework, there are two key problems. First, whether it is easy to model sensor network applications using the framework. Second, how to design efficient and accurate identification functions.

2.4 The Applicability of the Framework

The above framework is application independent, and thus can be used to model a large range of sensor networks. In this section, we use two examples to show its applicability.

Example 1: Localization in sensor networks. The approach for detecting compromised beacon nodes [Liu et al. 2005] described in section 2.1 can be modeled by our framework as follows.

- (1) **Observability graph.** The vertices of the observability graph only includes beacon nodes. There is a bi-directional edge between two beacon nodes, if they are close enough to probe each other. (Note that a detecting beacon node has to use a pseudoID to prevent a compromised beacon node from recognizing a probing query) Figure 2 shows the observability graph corresponding to Figure 1.
- (2) **Alerts.** If at time t a beacon node b_1 detects a bogus location claimed by a nearby beacon b_2 , it will send an alert (t, b_1, b_2) to the base station.

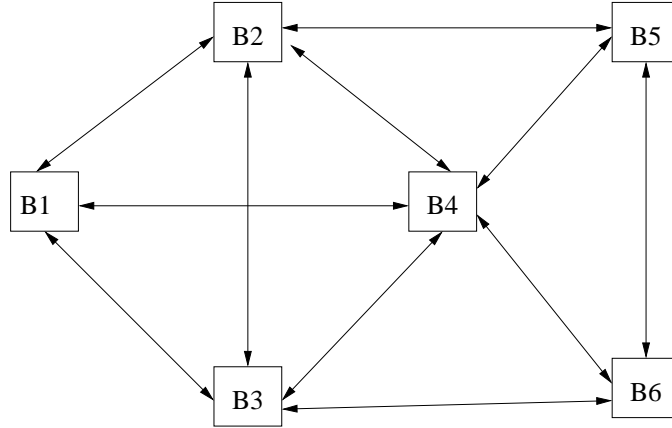


Fig. 2. The observability graph of sensor network localization

- (3) Sensor behavior model. The probability that an uncompromised beacon claims an incorrect location is determined by the resolution and the accuracy of the localization device. Let (x, y) be a beacon's actual location and (x', y') be its measured location. Assume $x' - x$ and $y' - y$ follow normal distribution $N(0, \sigma^2)$. Then the distribution of the distance d between (x, y) and (x', y') can be computed from bivariate transformation as $\frac{d}{2\pi^2\sigma^2} e^{-d^2/2\sigma^2}$. Suppose a beacon node is considered defective if d is larger than a pre-defined threshold ϵ . Then the reliability of beacon nodes is $r_m = P(d < \epsilon)$.
- (4) Observer model. Since a probe is always initiated by an observer, the observability rate of beacon nodes is $r_b = 1$. A beacon node's positive and negative detection rates can be derived from the given parameters of the detector. Specifically, the negative detection rate is the possibility when a benign beacon nodes reports correct location information, and its observer does not raise any alert. Liu et al. [Liu et al. 2005] proposed the false positive rate P_{fp} as a parameter of the detector, which is the rate that a benign beacon node raises alerts against another benign beacon node. So we have the negative detection rate $r_n = 1 - P_{fp}$. The positive detection rate is the probability that the observer successfully detects it when a malicious beacon nodes provides false location data. This is the same as the detection rate P_r as defined by Liu et al. in [Liu et al. 2005].

Since each observer probes a beacon node independently, one potential risk is that the beacon node may behave different to different observer, so that the number of alerts from different observers will be significantly different statically. This will make the first assumption in section 2.2 not valid anymore.

As mentioned above, the scheme of Liu et al. requires a detecting beacon node to use a completely new pseudoID for each probe. Therefore, a malicious beacon node cannot distinguish probes from different observers. Therefore, the number of alerts from uncompromised observers should still be statistically consistent during a period of time.

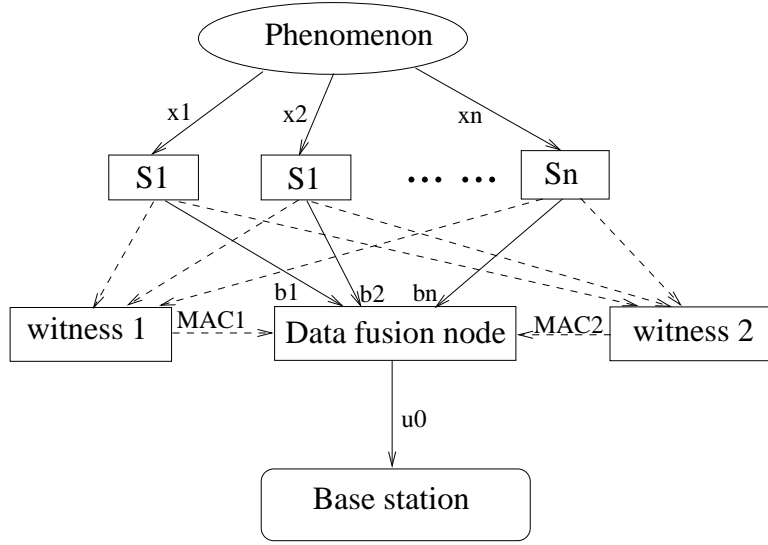
We intentionally omit the discussion of identification functions when modeling the above application with the proposed framework, since their work cannot handle false alerts from compromised nodes. Though inconsistencies between alerts may be discovered, no solu-

tion is provided to reason which nodes are compromised. We will show how to design an efficient identification function based on the above framework in the next section. Before doing so, we use another example to further demonstrate how to map a real sensor network application to the proposed framework.

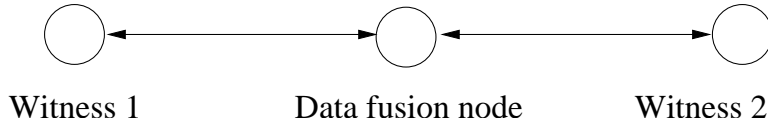
Example 2: Sensor network data fusion. Data fusion is a common technique to reduce communication in sensor networks. Instead of having all the raw data sent to the base station, they are sent to nearby data fusion nodes, which aggregate them and forward the result further to the base station. Compromised data fusion nodes may send bogus aggregated information to the base station. Du et al. [Du et al. 2003a] proposed a witness-based mechanism to detect compromised data fusion nodes. The basic architecture of their approach is shown in figure 3(a). For a set of sensor nodes, there are one dedicated data fusion node and several witness nodes, all of which can receive sensing data from sensor nodes in the set. A witness node performs the same aggregation as the data fusion node, but only forwards the MAC of the result (using its shared key with the base station) to the data fusion node. The data fusion node then sends its result along with the witnesses' MACs to the base station. The base station checks the consistency between the data fusion node's result and the MACs of witnesses. This scheme assumes a shared secret key between each witness and the base station. So only the base station can check whether a MAC agrees with the data fusion result. Clearly, if a witness is compromised, it may also send bogus MACs.

The above scheme can be modeled by our framework as follows.

- (1) Figure 3(b) shows the observability graph corresponding to figure 3(a). If there is a disagreement between the data fusion node and a witness, it is equivalent that they raise alerts against each other. In [Du et al. 2003a], sensor nodes collecting raw data are assumed to be trusted. Therefore, the observability graph only contains data fusion nodes and witnesses. The observability graph will be more complicated if we consider the situation where a node may serve as a witness or an aggregator for more than one set of sensors.
- (2) Alerts. If the MAC of a witness w_1 does not agree with the reported aggregation result of the data fusion node d_f for an event happened at time t , the base station can derive two alerts (t, w_1, d_f) and (t, d_f, w_1) .
- (3) Sensor behavior model. We only need to derive the reliability rate for the data fusion node and witnesses. Suppose the communication channels between sensor nodes are reliable. Then the data fusion node and witness nodes will always receive the same data from the sensor node and get the same aggregation results. Thus the reliability r_m of the data fusion node and the witness nodes is 1.
- (4) Observer model. If we assume the communication channels between sensor nodes are reliable, then the data fusion node and the witness nodes will always send data to the base station when an event happens. Thus the observability rate $r_b = 1$. The positive accuracy is the probability that the MAC sent by an uncompromised witness does not agree with the bogus aggregation result from a compromised data fusion node, which equals to $1 - 2^{-k}$, where k is the length of the MAC. The negative detection rate is the probability that the MAC from an uncompromised witness matches with the aggregation result from an uncompromised data fusion node, which equals to r_m .



(a) The architecture of the detection mechanism for data fusion



(b) The observability graph of the detection mechanism for data fusion

Fig. 3. Observability graphs for the data fusion example

Similar to the previous example we omit the discussion of identification functions here, as their work cannot identify the compromised nodes either. In the following, we are going to present our approach to identifying compromised sensor nodes based on the above framework. Generally, we first derive the expected alert pattern from the sensor behavior model and observer model. Then by comparing with the actual set of alerts, we identify those alerts which deviate from the expected behavior as abnormal. Finally combining this information with the observability relationship we design efficient identification functions such that the largest number of compromised nodes will be identified.

3. IDENTIFICATION OF COMPROMISED NODES

In this section, we present our approach to identifying compromised sensor nodes based on the above framework.

3.1 Overview of System Architecture

Let s_i be an observer of s_j . For each event at s_j , s_i will determine whether s_j has behaved correctly. If s_i believes s_j 's behavior is suspicious, then an alert will be raised and send back the base station. Here an event at s_j can be any behavior that the underlying detection system is interested in, according to different applications. For example, it can be a temperature value that s_j reports about the environment, or the fact that s_j has routed a packets

for others, or the location information from s_j in response to the disguised beacon node s_i . Note that the functionalities of event monitoring, alert generation and transmitting back to the based station are executed by the underlying detection mechanism, not by our model. Thus no additional communication and computation costs are imposed on the network.

In an ideal network, where the detection mechanism is completely accurate, and the sensing and observing capabilities of sensor nodes are perfect, if s_i raises an alert against s_j , then at least one of them is compromised. Otherwise, if both behave normally, there should be no alerts between them. However, sensors are not assumed to have perfect sensing and monitoring capabilities. Therefore, the base station cannot draw any definite conclusion from a single alert. Instead, it needs to observe the alert pattern during a certain period of time to discover suspicious activities with high confidence. Due to this reason, the base station breaks the total operation time of the network into time intervals (each is called a time window), and it reviews and reasons the alerts issued within each time window.

Within a given time window, the number of events of s_j is a random variable x , with a distribution $f_j(x)$. Given the sensor behavior model of s_j , observer model of s_i , and $f_j(x)$, we can derive the expected number of alerts raised by s_i against s_j in the time window, when both of them are uncompromised. Then the base station can compare the number of alerts actually raised by s_i against s_j with the expected number within this window. Only when the former is higher than the latter with statistical significance, should the base station consider it as abnormal.

For each pair of nodes that are considered abnormal, the base station will record an edge correspondingly in the observability graph (called an abnormal edge). As the base station has the knowledge of the global sensor topology, it can combine this information with all abnormal edges at hand to do further inference. Finally, we propose effective algorithms to reason the inferred graph and identify the compromised nodes.

Next, we will describe how the system works in details.

3.2 Alert Aggregation

Suppose s_j is a good node that functions normally according to the sensor behavior model, and s_i is a good observer to s_j , thus its behavior follows the observer model. For every event at s_j , with probability r_m it will behave correctly. The observer s_i can detect s_j 's behavior with probability r_b , and the probability that it issues a false alert stating that s_j behaves inappropriately is $1-r_n$, according to the observer model. With probability $1-r_m$, s_j will make mistakes, and s_i can detect this and issue a correct alert with probability r_p . These are the only two possibilities that s_i will raise an alert against s_j on a single event. Overall, the probability that s_i raises an alert against s_j on this single event is given by $C = r_b \cdot ((r_m \cdot (1 - r_n) + (1 - r_m) \cdot r_p))$. We interpret this as for each single event at s_j , there will be C alerts associated with it from s_i to s_j . In the sensor network localization, for example, each event is a probe from one beacon node to another beacon node. As shown in section 2.4, $r_b = 1$, $r_m = P(d < \epsilon)$, $r_n = 1 - P_{fp}$, and $r_p = P_r$. So $C = P(d < \epsilon) \cdot P_{fp} + (1 - P(d < \epsilon)) \cdot P_r$.

Let x denote the number of events at s_j within a time window T , which is a random variable, and $f_j(x)$ be the distribution of x . Then the distribution of alerts along the edge (s_i, s_j) (i.e., those raised by s_i against s_j) will be $f_{ij}(x) = f_j(\frac{x}{C})$. The expected number of alerts along (s_i, s_j) during the same period t will be $R_{ij}(t) = C \int^t f_{ij}(x)$.

During the time window T , if the number of alerts along the edge (s_i, s_j) is over $R_{ij}(T) + \delta$, then we say the edge (s_i, s_j) is an *abnormal edge* in the observability graph. Otherwise, (s_i, s_j) is *normal*. An abnormal edge can be interpreted as a definite claim from s_i that s_j is compromised. Similarly, a normal edge represents s_i 's endorsement that s_j is not compromised. Here the parameter $\delta > 0$. Given the distribution of the expected number of alerts during T , δ can be computed as the $x\%$ confidence that the number of alters are considered to be excessive. We leave the decision of $x\%$ to applications, as some applications may require a conservative reasoning, while others may not. The base station can even assign different $x\%$ to different observers according to their location, environment, functionality, etc. We will not elaborate the details since they are highly application specific.

Our framework does not require that the number of events of s_j monitored by each observer s_i to be the same. This will accommodate the underlying detector model to the largest extent. For example, in the beacon node application [Liu et al. 2005], different beacon node may probe a target beacon node different times within each time window. So the event distribution of s_j observable to s_i is different for each node s_i . But we can still derive the expected number of alerts along each edge (s_i, s_j) , and tell if the edge is abnormal or not using the above reasoning.

Note that it is possible that s_j is not compromised but malfunctioning. But in this case, we treat s_j as compromised anyway since information or services from s_j cannot be trusted anymore.

3.3 Graph Inference

Given a set of abnormal and normal edges, many existing trust functions in the literature can be applied to infer each sensor node's trustworthiness. However, since those functions are designed for general decentralized systems, they do not take into consideration the interaction topology between sensor nodes, which in fact provides valuable information to directly identify compromised sensor nodes. For example, suppose that there are no more than k compromised sensor nodes in a sensor network. If more than k abnormal edges involve a sensor node s , then we can know for sure that s is compromised. How to take advantage of such information to more effectively identify compromised sensor nodes is the focus of the rest of the paper.

Given an abnormal edge (s_i, s_j) , either s_i is compromised and raises many bogus alerts against s_j , or s_j is compromised and its malicious activities are observed by s_i , or both. Otherwise, the edge should be normal. Further suppose there is an additional normal edge (s_l, s_j) . Then one of s_i and s_l must be compromised. Otherwise, the two edges should be consistent with each other since s_l and s_i observe the activities of the same node s_j during the same period of time, according to our assumptions in section 2.2. Information of such inconsistency is essential to identify compromised nodes.

Definition 3.1. Given an observability graph $G(V, E)$, let E_a and E_n be the set of abnormal edges and normal edges in G respectively. We say that two sensor nodes s_i and s_j form a *suspicious pair* if one of the following holds:

- (1) $(s_i, s_j) \in E_a$ or $(s_j, s_i) \in E_a$;
- (2) There exists a sensor node s' , such that either $(s_i, s') \in E_a$ and $(s_j, s') \in E_n$, or $(s_i, s') \in E_n$ and $(s_j, s') \in E_a$.

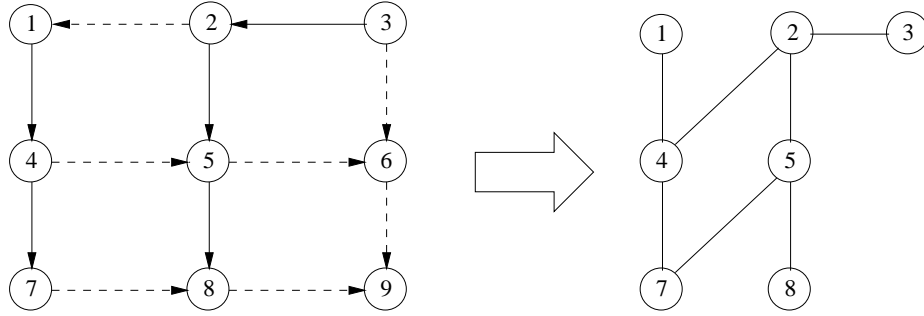


Fig. 4. An observability graph and its corresponding inferred graph

Let $\{s_1, s'_1\}, \dots, \{s_k, s'_k\}$ be the suspicious pairs derived from an observability graph G . The *inferred graph* of G is an undirected graph $I(V', E')$ such that $V' = \bigcup_{1 \leq i \leq k} \{s_i, s'_i\}$ and $E' = \{\{s_i, s'_i\} \mid 1 \leq i \leq k\}$.

Intuitively, if (s_i, s_j) are a suspicious pair, then at least one of them is compromised. Note that if a pair of node is not suspicious, it does not mean that they are both uncompromised. It only means we cannot infer anything about them.

The left part of figure 4 shows an observability graph, where abnormal and normal edges are represented by solid and dashed edges respectively. The corresponding inferred graph is shown in the right part of the figure. Note that an inferred graph may not be connected.

From the figure one may wonder that, since $(2, 1)$ is normal and $(2, 5)$ is abnormal, shouldn't $\{1, 5\}$ be a suspicious pair too? The answer is no since it is possible that node 2 is compromised, and selectively issues bogus alerts against node 1 but not node 5, even though both node 1 and 5 are uncompromised. Similarly, a compromised node may sense data normally but issue bogus alerts or vice versa. So $\{1, 3\}$ is not a suspicious pair even though $(3, 2)$ is abnormal and $(2, 1)$ is normal. In other words, transitivity does not hold when constructing suspicious pairs.

Definition 3.1 in fact offers two inference rules to identify pairs of sensor nodes such that at least one of them is compromised. A natural question is whether the two inference rules are *complete*. In other words, are there any other suspicious pairs existing that cannot be identified using the above two inference rules? The answer is no unless we impose further assumptions regarding the behavior of a compromised node.

A node s takes multiple responsibilities in a sensor network. Besides acting as a common sensor node, it also acts as observers of several other nodes s_1, \dots, s_k . Without any assumptions on the application of the sensor network, a compromised nodes may behave arbitrarily when fulfilling different responsibilities. As the above example shows, a node may report false sensing data and meanwhile act normally as an observer, or vice versa. It may also issue bogus alerts to one node but not to others. Thus, the behavior of s in a time window can be represented as a vector $behave(s) = (sensing(s), obs(s, s_1), \dots, obs(s, s_k))$. Each element in the vector can either take the value g or b , corresponding to normal and abnormal behavior respectively. A behavior of s is an assignment of values to the vector of s . Note that, since a compromised node may behave arbitrarily, there is no correlation between the value of each element. They can be assigned independently. Clearly, if at least one of the elements in the vector must be assigned to b , then the node should be

considered as compromised. Otherwise, the node is considered as uncompromised, as long as there exist assignments that all its elements can be set to g .

Let $G(V, E)$ be an observability graph marked with normal and abnormal edges, and $\mathcal{B} = \{behave(s_1), \dots, behave(s_n)\}$ be a set of behaviors of all the nodes in G . We say \mathcal{B} is *consistent* with G if the two conditions in definition 3.1 are satisfied. That is, if there is an abnormal edge (s_1, s_2) , then either $sensing(s_2) = b$ or $obs(s_1, s_2) = b$. Further, if (s_1, s) is normal and (s_2, s) is abnormal, then either $obs(s_1, s) = b$ or $obs(s_2, s) = b$.

THEOREM 3.2. *If two nodes s_1 and s_2 are not identified as a suspicious pair by definition 3.1, then there exists a consistent set of behaviors of all the nodes where both s_1 and s_2 are uncompromised.*

PROOF. We first construct a consistent behavior set \mathcal{B} as follows. For each node s , we set $sensing(s) = g$. For each abnormal edge (s', s) , let $obs(s', s) = b$. For each normal edge (s'', s) , let $obs(s'', s) = g$. It is easy to see that the resulting behavior set is consistent.

If in \mathcal{B} both s_1 and s_2 are uncompromised, that is, all elements in $behave(s_1)$ and $behave(s_2)$ are assigned g , then we are done. Otherwise, s_1 must have an element $obs(s_1, s')$ that is assigned to b . Consider each node s' such that $obs(s_1, s') = b$. Because s_1 and s_2 are not identified as a suspicious pair by definition 3.1, $s' \neq s_2$. If s_2 is also an observer of s' , we must have $obs(s_2, s') = b$ as well. Therefore, we can set $sensing(s')$ to be b , and $obs(s_1, s')$ to be g . We also set $obs(s_2, s')$ to be g when s_2 is also an observer of s' . Next for each edge (s'', s') , if it is abnormal, we set $obs(s'', s) = g$. Otherwise, set $obs(s'', s) = b$. By doing so, we remove one bad assignment in the behavior vector of s_1 , and will not change any good assignment in the behavior vector of s_2 to bad. And the resulting behavior set is still consistent.

We keep doing the above modification for each node that s_1 and s_2 observe. The final behavior set is consistent and both s_1 and s_2 are uncompromised. \square

Note that for a specific sensor network application, it may have further constraints on the behavior of a compromised node. For example, a node's observing behavior may be tied to its sensing behavior, i.e., they have to be normal or abnormal at the same time. Without such application specific properties, we may identify more suspicious pairs, which is outside of the scope of this paper, as we focus on a general application-independent framework for identifying compromised nodes.

Next, we discuss how to identify compromised nodes when taking the security estimation into consideration.

Clearly, if a sensor node does not appear in the inferred graph, then its behavior is consistent with the sensor behavior model and observer model, and thus should be considered uncompromised. Hence, we concentrate on identifying compromised sensor nodes among those involved in the inferred graph.

Definition 3.3. Given an inferred graph $I(V, E)$ and a security estimation K , a *valid assignment* with regard to I and K is a pair (S_g, S_b) , where S_g and S_b are two sets of sensor nodes that satisfies all of the following conditions:

- (1) S_g and S_b is a partition of V , i.e., $S_g \cup S_b = V$ and $S_g \cap S_b = \emptyset$;
- (2) For any two sensor nodes s_i and s_j , if $s_i \in S_g$ and $s_j \in S_g$, then $\{s_i, s_j\} \notin E$; and
- (3) $|S_b| \leq K$.

Intuitively, a valid assignment corresponds to one possible way that sensor nodes are compromised, that is, when they raise false alerts and conduct abnormal activities, the resulting inferred graph is I . S_g and S_b contains the uncompromised and compromised nodes respectively. For a given inferred graph and a security estimation K , there may exist many valid assignments. Obviously the common nodes in all possible assignments are always compromised, and others may or may not be compromised, depending on which assignment is true for the actual system. This inspires us that an optimal algorithm is to identify the common nodes in all possible assignments, thus it will identify the largest number of truly compromised nodes, and does not introduce any false positives.

Definition 3.4. Given an inferred graph $I(V, E)$ and a security estimation K , let $\{(S_{g1}, S_{b1}), \dots, (S_{gn}, S_{bn})\}$ be the set of all the valid assignments with regard to I and K . We call $\bigcap_{1 \leq i \leq n} S_{bi}$ the *compromised core* of the inferred graph I with security estimation K , denoted $CompromisedCore(I, K)$. Similarly, $\bigcap_{1 \leq i \leq n} S_{gi}$ is called the *uncompromised core* of I with security estimation K , denoted $UncompromisedCore(I, K)$.

Definition 3.5. Let I be the inferred graph, given an observability graph G , a sensor behavior model, an observer model, a security estimation K , and a set of alerts during a time period T . We say an identification function is F *optimal* if and only if F always returns $CompromisedCore(I, K)$.

If an identification function is optimal, then it identifies the largest number of compromised nodes without introducing any false positives. In other words, if a function F' returns any node s not in the compromised core, then we can always find a valid assignment such that s is not compromised in the assignment. This implies F' may introduce false positives. Given the general framework, *one key problem is thus to develop algorithms that efficiently compute $CompromisedCore(I, K)$.*

On the other hand, though introducing no false positives, the compromised core may not achieve high detection rates since there may exist suspicious pairs whose nodes are not included in the compromised core. Thus, *another key problem is to seek techniques that further eliminate compromised nodes without causing many false positives.*

In summary, our approach is composed of two phases. In the first phase, we compute or approximate the compromised core, identifying those nodes that are definitely compromised. In the second phase, we tradeoff accuracy for eliminating more compromised nodes.

3.4 The Algorithm to Identify Compromised Sensor Nodes

Though collusion between compromised nodes is good for an attacker, an identification function should not rely on any assumptions of collusion models. Otherwise, an attacker may easily defeat the identification algorithm by slightly changing the behavior of compromised nodes and making the collusion assumption invalid. For example, even if s_1 issues a lot of alerts against s_2 , we cannot conclude that one of them is compromised and the other is not. It is possible that both of them are compromised and the attacker just wants to confuse the identification function.

On the other hand, no matter how compromised nodes collude, it always holds that a suspicious pair contains at least one compromised node. This property helps us derive the lower bound of the number of compromised nodes.

LEMMA 3.6. *Given an inferred graph $I(V, E)$, let V_I be a minimum vertex cover of I . Then the number of compromised nodes is no less than $|V_I|$.*

PROOF. When we assign the nodes into S_g, S_b , for each edge E_{ij} , at least one end point of them is in S_b . So S_b is essentially a vertex cover for I . $|S| \geq |V_I|$. \square

We denote the size of the minimum vertex covers of an undirected graph G as C_G . Given a sensor node s , the neighbors of s in an inferred graph I is denoted \mathcal{N}_s . Further, let I'_s denote the graph after removing s and its neighbors from I . We have the following theorem for identifying compromised sensor nodes.

THEOREM 3.7. *Given an inferred graph I and a security estimation K , for any node s in I , $s \in \text{CompromisedCore}(I, K)$ if and only if $|\mathcal{N}_s| + C_{I'_s} > K$.*

PROOF. \Rightarrow : Suppose there exist some s that satisfies $|\mathcal{N}_s| + C_{I'_s} > K$, and $s \in S_g$ for some assignment (S_g, S_b) . Then we will have $N_s \subseteq S_b$. According to Lemma 3.6, we know the minimum number of malicious nodes in I'_s is $C_{I'_s}$. So we have $|S_b| \geq |\mathcal{N}_s| + C_{I'_s} > K$. This contradicts with constraint 3 of Definition 3.3. So s must be in S_b under any assignment. Thus $s \in \text{CompromisedCore}(G, K)$.

\Leftarrow : If there exists $s \in \text{CompromisedCore}(G, K)$ that satisfies $|\mathcal{N}_s| + C_{I'_s} \leq K$. Then we can always construct an assignment of S_b : $S_b = \mathcal{N}_s \cup V_{I'_s}$. This assignment will satisfy all constraints of Definition 3.3, but it also introduces a contradiction: $s \in \text{CompromisedCore}(G, K)$, and $s \notin S_b$. So if $s \in \text{CompromisedCore}(G, K)$, it must satisfy $|\mathcal{N}_s| + C_{I'_s} > K$. \square

Intuitively, if we assume a sensor node s is uncompromised, then all its neighbors in I must be compromised. According to lemma 3.6, there are at least $|\mathcal{N}_s| + C_{I'_s}$ compromised nodes, which should be no more than the security estimation K . Otherwise, s must be compromised. Meanwhile, if $|\mathcal{N}_s| + C_{I'_s} \leq K$, we can always construct a valid assignment for I with regard to K , where s is assigned as uncompromised, which means s is not in $\text{CompromisedCore}(I, K)$.

By theorem 3.7, the algorithm to identify $\text{CompromisedCore}(I, K)$ is straightforward. For each node s , we check whether $|\mathcal{N}_s| + C_{I'_s}$ is larger than K . Unfortunately, this algorithm is in general not efficient since the minimum vertex covering problem is NP-complete. In theory we also have to compute the minimum vertex cover of a different graph when checking each node.

Thus, we seek efficient algorithms to approximate the size of minimum vertex covers. To prevent false positives, we are interested in deriving a good *lower bound* of the size of minimum vertex covers, a goal different from that of many existing approximation algorithms. In this paper, we choose the size of maximum matchings of I as such an approximate. We denote the size of the maximum matchings of an undirected graph G as M_G .

LEMMA 3.8. *Given an undirected graph G , $M_G \leq C_G \leq 2M_G$. And the bounds are tight, [Vazirani 2001].*

COROLLARY 3.9. *Given an inferred graph I and a security estimation K , for any node s in I , if $|\mathcal{N}_s| + M_{I'_s} > K$, then $s \in \text{CompromisedCore}(I, K)$. \square*

A maximum matching of an undirected graph can be computed in polynomial time [Micali and Vazirani 1980]. Algorithm 1 shows an efficient algorithm to approximate the compromised core. Since this algorithm does not assume any specific collusion model among compromised nodes, we call it the *general identification algorithm*.

Algorithm 1 AppCompromisedCore(I, K)

```

//Input:  $I$  is an inferred graph
//       $K$  is a security estimation
//Output: the compromised core of  $I$  with  $K$ 
 $S_b = \emptyset$ 
For each sensor node  $s$  in  $I$ 
    Let  $n_s$  be the number of neighbors of  $s$ 
    Let  $m = M_{I'_s}$ 
    If  $n_s + m > K$ 
         $S_b = S_b \cup \{s\}$ 
Return  $S_b$ 

```

THEOREM 3.10. *The complexity of the algorithm AppCompromisedCore is $O(mn\sqrt{n})$, where m is the number of edges and n is the number of vertices in an inferred graph.*

PROOF. From [Micali and Vazirani 1980], the complexity of finding the maximum matching of I'_s of any node s is $O(m\sqrt{n})$. Algorithm AppCompromisedCore will go through each node s in the system, so the total complexity is $O(mn\sqrt{n})$. \square

3.5 Further Elimination of Compromised Sensor Nodes

The above algorithm does not introduce any false positives. Compromised nodes identified by the above algorithms may be safely excluded from the networks through, e.g., key revocation [Du et al. 2003b; Liu and Ning 2003]. However, there may still be suspicious pairs left that do not include any nodes in the compromised core. We call the graph composed of such pairs the *residual graph*.

We may tradeoff accuracy for eliminating more compromised nodes. Since a suspicious pair contains at least one compromised node, identifying both nodes as compromised will introduce at most one false positive. By computing the maximum matching of a residual graph and treating them as compromised, the false positive rate is bounded by 0.5. Note that this is the best we can do based on the information provided by the general framework. In order to reduce this worst-case false positive rate, application specific information or assumptions are needed.

The complexity of this phase is bounded by $O(m\sqrt{n})$ [Micali and Vazirani 1980], where m and n are the numbers of edges and vertices in an inferred graph respectively.

In summary, given an inferred graph and a security estimation, our approach is first to approximate its compromised core. We then compute the maximum matching of the residual graph, and further eliminate compromised nodes.

3.6 Identification of Colluding Compromised Sensor Nodes

In general, collusion among compromised nodes will enable attackers to have a stronger influence on the sensor network, and make it harder for the base station to identify them. The general identification algorithm does not assume any knowledge of the collusion among compromised nodes. Therefore, the worst case would be that the compromised nodes do collude, but the base station has to be conservative and use the general identification algorithm.

In this section, we consider the situation where the base station is aware of how compro-

mised nodes will collude. Note that we are not advocating that this is a realistic situation, as in practice an attacker's strategy is most likely unknown. Instead, we intend to investigate how much the base station can do better with this extra knowledge, which will help us better understand the identification capability of the general algorithm. In other words, the identification capability of the base station with the extra knowledge of the attacker's strategy gives us an upper bound of that of the general identification algorithm.

For simplicity, in this section, we assume a strong collusion model where compromised nodes are coordinated with each other when monitoring events and raising alerts. In other words, compromised nodes do not raise alerts against each other, and alerts against a node from two compromised nodes are consistent. It is easy to see that in this case the inferred graph is bipartite, since an edge can only be between an uncompromised node and a compromised one. We call inferred graphs in this situation *collusion-inferred graphs*.

Let G_1, \dots, G_k be the connected components of a collusion inferred graph I . For each G_i , $1 \leq i \leq k$, let A_i, B_i denote the two sets of disjoint vertices in G_i such that vertices in one set are adjacent only to vertices in the other set. We call (A_i, B_i) the *two-color partition* of G_i . Without loss of generality, we assume $|A_i| \geq |B_i|$. Clearly, since I is collusion inferred, vertices in the same set are either all compromised or all uncompromised. Further, if vertices in one set are compromised (uncompromised), then vertices in the other set are uncompromised (compromised). Therefore, $\sum_{1 \leq i \leq k} |B_i|$ becomes a lower bound of the number of compromised nodes in a network.

LEMMA 3.11. *Let $I(V, E)$ be a collusion-inferred graph and K be a security estimation. Given the two-color partition (A_i, B_i) of any connected component of I , $1 \leq i \leq k$, let Q_i be either A_i or B_i . We have $Q_i \subseteq \text{UncompromisedCore}(I, K)$ if and only if $|Q_i| + \sum_{j \neq i} \min(|A_j|, |B_j|) > K$.*

PROOF. \Rightarrow : For simplicity we define the notation of sets in such a way that $|A_j| \geq |B_j|$ for all j as before. Suppose there exists $Q_i \subseteq \text{UncompromisedCore}(G, K)$, which satisfies $|Q_i| + \sum_{j \neq i} |B_j| \leq K$. Then from definition 3.4, $Q_i \in S_{gr}$ for all $1 \leq r \leq n$, where n is the number of valid assignments. Suppose the corresponding set in the same coloring partition with Q_i is P_i . Since P_i and Q_i are bipartite, we know that $P_i \in S_{br}$ for all $1 \leq r \leq n$, thus $P_i \subseteq \text{CompromisedCore}(G, K)$. Now we look at a particular assignment (S_{gx}, S_{bx}) , $1 \leq x \leq n$, such that S_{bx} contains Q_i and all B_j , $1 \leq j \leq k$ except Q_i , where k is the number of two-color partitions, and S_{gx} contains P_i and all A_j , except Q_i . This assignment satisfies all constraints of definition 3.3, but it also introduces contradiction: $Q_i \subseteq \text{UncompromisedCore}(G, K)$, thus $Q_i \subseteq S_{gj}$ for all $1 \leq j \leq n$, but $Q_i \not\subseteq S_{gx}$ in this assignment. So if we have $Q_i \subseteq \text{UncompromisedCore}(G, K)$, it must satisfy $|Q_i| + \sum_{j \neq i} |B_j| > K$.

\Leftarrow : Suppose there exists some Q_i that satisfies $|Q_i| + \sum_{j \neq i} |B_j| > K$, and Q_i can be in some S_{by} , $1 \leq y \leq n$. Since each coloring partition will have a set in S_{by} , S_{by} will be composed of the following sets: $S_{by} = \{Q_i, A_r, \dots, A_s, B_t, \dots, B_m\}$. We thus have $|S_{by}| > |Q_i| + \sum_{j \neq i} |B_j| > K$, since $|A_i| \geq |B_i|$ for $1 \leq i \leq k$. This contradicts with constraint 3 in Definition 3.3. So Q_i must be in S_{gy} for all $1 \leq y \leq n$. Thus $Q_i \subseteq \text{UncompromisedCore}(G, K)$. \square

Intuitively, if Q_i is compromised, then $|Q_i| + \sum_{j \neq i} |B_i|$ gives a lower bound of the number of compromised nodes, which should be no more than K .

COROLLARY 3.12. *Given any connected component G_i of a collusion inferred graph I , and a security estimation K , $B_i \not\subseteq \text{UncompromisedCore}(I, K)$.*

PROOF. Suppose there exists some $B_i \subseteq \text{UncompromisedCore}(I, K)$, then we know $B_i \in S_{g_i}$ for all assignment $1 \leq i \leq n$. Thus $A_i \in S_{b_i}$ for all assignment $1 \leq i \leq n$, so we have $A_i \subseteq \text{CompromisedCore}(I, K)$. Now from Lemma 3.11, we have $|B_i| + \sum_{j \neq i} |B_j| > K$. Since $|A_i| \geq |B_i|$, we also have $A_i + \sum_{j \neq i} |B_j| > K$. Again from Lemma 3.11, A_i must be in $\text{UncompromisedCore}(I, K)$. This introduces contradiction. So $B_i \not\subseteq \text{UncompromisedCore}(I, K)$. \square

THEOREM 3.13. *Given a collusion inferred graph I with k connected components, and a security estimation K , let (A_i, B_i) , $1 \leq i \leq k$, be the two-color partition of each connected component. Let $S_b = \{B_i \mid |A_i| + \sum_{j \neq i} |B_j| > K\}$, and $S_a = \{A_i \mid |A_i| + \sum_{j \neq i} |B_j| > K\}$. Then $\text{CompromisedCore}(I, K) = \bigcup_{B_i \in S_b} B_i$ and $\text{UncompromisedCore}(I, K) = \bigcup_{A_i \in S_a} A_i$.*

PROOF. This is the direct observation based on Lemma 3.11 and Corollary 3.12. \square

Given theorem 3.13, it is straightforward to develop the algorithm that efficiently compute $\text{CompromisedCore}(I, K)$. The pseudocode is shown as Algorithm 2.

Algorithm 2 CollusionCompromisedCore(I, K)

```
//Input:  $I$  is a collusion inferred graph
//       $K$  is a security estimation
//Output: the compromised core of  $I$  with  $K$ 
 $S_b = \emptyset$ 
For each connected component  $G_i$  of  $I$ 
    Let  $(a_i, b_i)$  be the two-color partition of  $G_i$ 
    If  $|a_i| + \sum_{j \neq i} |b_j| > K$ 
         $S_b = S_b \cup b_i$ 
Return  $S_b$ 
```

THEOREM 3.14. *The complexity of algorithm CollusionCompromisedCore is $O(n+m)$, where n and m are the numbers of vertices and edges respectively in an collusion inferred graph.*

PROOF. The complexity of finding all the two-color partitions from the collusion inferred graph is $O(m)$, since we need to go through each edge to construct the partitions. The complexity of identifying the compromised core is $O(n)$, since the number of the connected components is in the order of $O(n)$, and we need to go through each component. Thus the total complexity of algorithm CollusionCompromisedCore is $O(n+m)$. \square

After we have identified the compromised core for the bipartite graph, we achieve the goal of identifying the largest number of compromised nodes without introducing any false positives. For the residual graph, we can use the same maximum matching algorithm as we proposed in section 3.5 to tradeoff false positives for detection rates.

4. EXPERIMENTS

In this section, we design a set of experiments to evaluate the effectiveness of our algorithms.

4.1 Experiment Methodology

We simulate a sensor network deployed to monitor the temperature of an area of $100m \times 100m$. For simplicity, we assume sensor nodes are randomly distributed in the area. We adopt a simple detection mechanism. If the distance between two sensor nodes is within 10 meters, and the temperatures reported by them differ by more than $1^\circ C$, the base station infers that each of them raises an alert against the other. In other words, two nodes are observers of each other if they are within 10 meters. We assume that, once a network is deployed, sensors' location information can be collected through localization techniques. Therefore, the base station is able to construct an observability graph accordingly.

Sensor nodes report temperatures to the base station once per minute, and the sensed data follows a normal distribution $N(\mu, \sigma^2)$, where μ is the true temperature at a location, and $\sigma = 0.2$, which is consistent with the accuracy of typical thermistors in sensor network [Crossbow Technology Inc. 2003]. The overall operation time is 24 hours, but the base station will review the system every one hour. So within each time window of one hour, there will be 60 data reports from each sensor node. In the simulation, we do not consider such effects as the message loss, as it has already been modelled by the parameters as r_b , r_m and r_n .

Unless otherwise stated, we assume $10\% \sim 15\%$ of the nodes in the network are compromised. The security estimation is $K = 15\%N$, where N is the total number of nodes in the network. The goal of the attacker is to raise the temperature reading in the area. The attacker may choose to either randomly compromise nodes in the field, in which case no local majority is formed, or compromise nodes centered around a position (x, y) , following a normal distribution with a standard deviation σ_d . The latter corresponds to the case of local majority, and the parameter σ_d controls its strength. The smaller σ_d is, the closer compromised nodes are to each other, and thus the stronger the local majority is. We call σ_d the *concentration* of compromised nodes.

The evaluation metrics of the experiments are the detection rates and false positive rates of the proposed identification approach. In this section we show the detection rates and false positive rates separately in different figures. We do not plot the Receiver Operating Characteristic curve(ROC) of our algorithm since given an observability graph and a set of alerts, there is no other parameters in our approach to adjust the identification results.

The effectiveness of our algorithm will be affected by the actual deployment of the sensor nodes, especially, the observability graph. Thus for each simulation, we generate 10 random deployment of sensor networks, following the parameters specified for each one's setup. The detection rates and false positive rates are averaged over the results on these 10 random networks. We have identified through simulations that the 95% confidence interval of the detection rate and false positive rates are within 1% across the 10 random networks, so we believe 10 should be enough.

We first act conservatively, and assume that compromised nodes in fact follow the strong collusion model, but the base station does not know this, and cannot utilize any knowledge of the collusion. Compromised nodes all report the same false temperature so that there are no alerts between them. We evaluate the effectiveness of the general *AppCompromised-*

Core algorithm followed by the maximum matching approach. We call it *general+mm* in the experiment. As a comparison, we will also see how well we can do if we know the fact that the attackers are colluding, and evaluate the *CollusionCompromisedCore* algorithm followed by the maximum matching approach. We call it *bipartite+mm* in the experiment.

We further compare our approach with the simple voting mechanism. Among a node s 's neighbors, if those that raise alerts against s are more than those that do not, then s is considered as compromised.

We also compare our approach with EigenTrust [Kamvar et al. 2003] and PeerTrust [Xiong and Liu 2002], two well-known reputation-based trust functions for P2P systems. Though there are many trust functions proposed for P2P networks and semantic webs, many of them are decentralized and subjective, in the sense that an entity's trust value varies, depending on who is the trust evaluator [Lee et al. 2003; Richardson et al. 2003; Yu and Singh 2002; Golbeck and Hendler 2004]. They are not suitable for centralized identification of compromised nodes in sensor networks. EigenTrust and PeerTrust both derive a unique global trust value for each entity from that entity's transaction records. Applied in sensor networks, the global trust values can be used to rank sensor nodes, where those with low trust values are identified as compromised. Therefore, they can be compared with *general+mm* and *bipartite+mm*.

The idea of EigenTrust is similar to the PageRank algorithm [Lawrence et al. 1998]. In EigenTrust, the number of satisfactory and unsatisfactory transactions between each pair of entities is collected to construct a matrix. One special property of the matrix is that the entries in each row add up to 1. The matrix is repetitively multiplied with an initial vector until it converges. The initial vector is a pre-defined parameter that corresponds to the default trust value of each entity. Each entry in the converged vector represents an entity's final trust value. Since the trust values of all nodes always add up to 1, $1/N$ is the average trust value of sensor nodes in the network. In the experiment, we identify the K nodes with the lowest trust values as compromised, unless their trust values are over $1/N$.

In PeerTrust, an entity's trust value is the normalized number of satisfactory services over the total number of transactions in which the entity takes part. It is further weighted by the trust values of the transaction feedback issuers and the quantity of each transaction. Similar to EigenTrust, the global trust values for all entities are also obtained through iterative matrix multiplication until it converges. In general, PeerTrust can be viewed as a majority voting scheme weighted by voters' trust values.

Traditional fault diagnosis PMC model are not applicable to sensor networks as we have analyzed in section 1, and we do not compare them.

We conduct four sets of experiments to evaluate the impacts of the following factors on the effectiveness of the above approaches.

4.2 Local Majority

Figure 5 shows the effectiveness of different approaches when compromised nodes form local majorities. The number of nodes in the network is set to be 200. The concentration of compromised nodes is varied from 5 to 100. When compromised nodes are extremely close to each other, they essentially form a cluster. Suspicious pairs only involve those compromised nodes near the edge of the cluster. Those near the center of the cluster do not appear in the inferred graph, and thus cannot be identified. We note that even in this situation the *general+mm* approach can still achieve detection rate over 0.6, mainly due to the maximum matching approach in the second phase. When the compromised

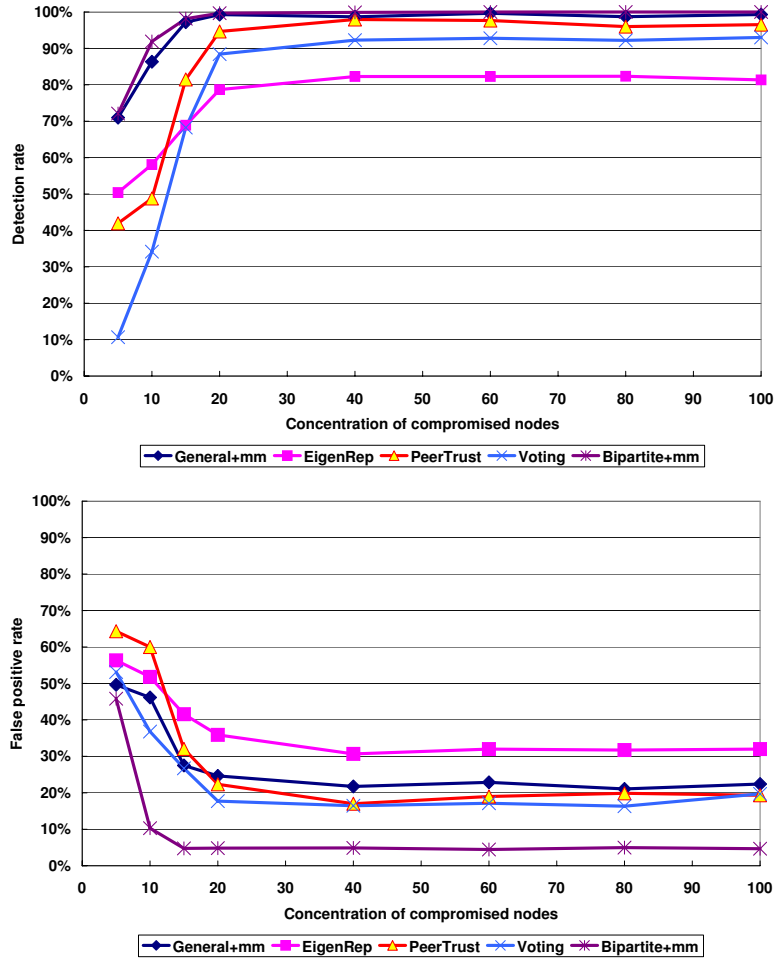


Fig. 5. The impact of the concentration of compromised sensor nodes

nodes are less concentrated, the general identification algorithm enables the base station to quickly identify almost all the compromised nodes. That is why we see a quick rise in general+mm’s detection rate and a sharp drop in its false positive rate.

When compared with other schemes, we have the following observations. First, EigenTrust seems to be inferior to general+mm, bipartite+mm and PeerTrust. The reason is that EigenTrust relies on the existence of pre-trusted peers to identify malicious collectives, which correspond to colluding compromised nodes in our setting. Without pre-trusted peers, it cannot significantly distinguish malicious entities from good ones. That is why we see an upper bound of the detection rate of EigenTrust even when compromised nodes do not form a strong local majority.

Second, we notice that when the concentration is over 20, PeerTrust and voting mechanism actually yield comparable detection rate to that of general+mm with a little bit lower false positive rates. A closer examination of the network reveals that, with 200 nodes in the network, the average number of observers for each node is around 3. When the concen-

tration is 20, among the neighbors of a compromised node, on the average no more than 1 neighbor is compromised. In other words, when the concentration is over 20, compromised nodes seldom form local majorities. In this case PeerTrust and simple voting, both relying on majority voting mechanisms, are more likely to assign low trust values to compromised nodes or label them as compromised nodes directly. For general+mm, each identified compromised nodes in the second phase will result in the sacrifice of an uncompromised nodes, resulting in higher false positive rates.

Third, when the compromised nodes form strong local majorities (i.e., the concentration is smaller than 20), general+mm yields much higher detection rates and lower false positive rates than PeerTrust. And the simple voting has the poorest detection rate as low as 10%, as it does not do any reasoning on the credibility of the feedback. This is an important advantage of our approach. In sensor networks, it is always cost-effective for attackers to compromise a small portion of the network, and make them collude. Otherwise either they have to compromise a large portion of the network, which is very costly and often not feasible, or they do not collude, in which case any voting-based algorithm can identify most of the compromised nodes, as shown above. So it is important that an identification algorithm performs well even when local majorities are formed by compromised nodes. From the experiment we see when collusion is the strongest, although the false positive rate of our algorithm is close to 50%, it is still the lowest among all solutions, and also achieves the highest detection rate.

Last, bipartite+mm always has the highest detection rate and lowest false positive rate. This shows when we have the exact information of the attackers behavior, the dedicated identification algorithm can be very effective.

4.3 Sensor Node Density

In the next experiment, we vary the number of sensor nodes in the area from 50 to 200. As we have seen from previous experiment, when there is no collusion among compromised nodes, all algorithms have high detection rate and false positive rate, and our algorithm outperforms other 3 algorithms when there exists a strong collusion. Thus we report the experimental results which set the concentration of compromised nodes to be 15, a not too strong case. We have also tried other concentration parameters and observed similar trends. Figure 6(a) and 6(b) show respectively the detection rate and the false positive rate of each approach.

We see that when the number of sensor nodes increases, all the approaches achieve better detection rates and less false positive rates. Intuitively, the more densely sensor nodes are deployed, the more observers each node has, and thus the more likely an abnormal activity is detected. The second observation is that general+mm and bipartite+mm do not achieve high detection rates when sensor nodes are deployed very loosely. This is because in this situation many nodes do not have any observers. There are too few alerts for the base station to identify compromised nodes definitely. A further study tell us that averagely when there are 100 nodes in this area, 16.6 nodes have only one observer. When there are 150 nodes in the area, 10 of them have only one observer. When we deploy 200 nodes in the area randomly, only 5.2 nodes averagely will have one observer.

The third observation, we see that general+mm detects much more compromised nodes with similar false positives than all other approaches, due to the fact that it takes the unique properties of sensor networks into consideration so that it is more resilient to compromised nodes forming local majorities. Finally, the bipartite+mm algorithm is still the best as it

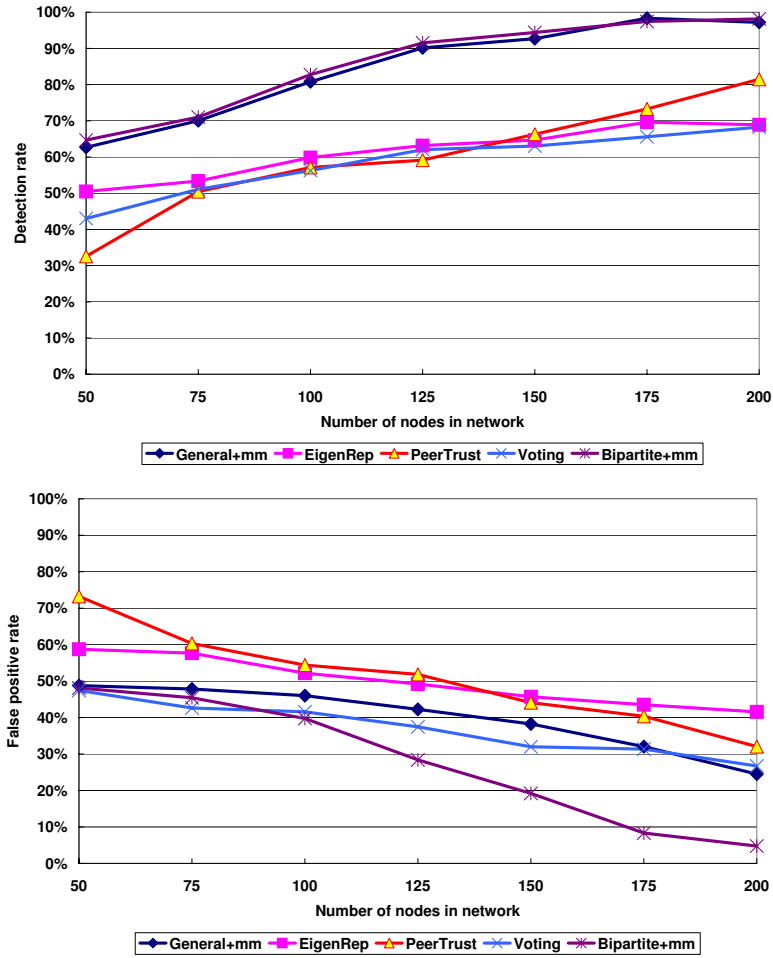


Fig. 6. The impact of the deployment density of sensor nodes

has additional information about attacker's behavior.

4.4 Accuracy of The Security Estimation

The security estimation gives an upper bound of the number of compromised nodes. It is an important parameter for identification functions to infer compromised nodes. An accurate security estimation is not expected to always be available. The next experiment evaluates how the accuracy of a security estimation affects the effectiveness of different approaches. The total number of sensor nodes are set to be 200. The number of compromised nodes is 20, and their concentration is set to be 15. The security estimation is varied from 20 to 40. Voting mechanism is not evaluated in this experiment as it does not involve this parameter when identifying compromised nodes. The experiment results are shown in figure 7.

We see that general+mm still achieve very high detection rates even when the accuracy of the security estimation decreases. When the security estimation is accurate, most of the compromised nodes are identified by the algorithm in the first phase, producing few false

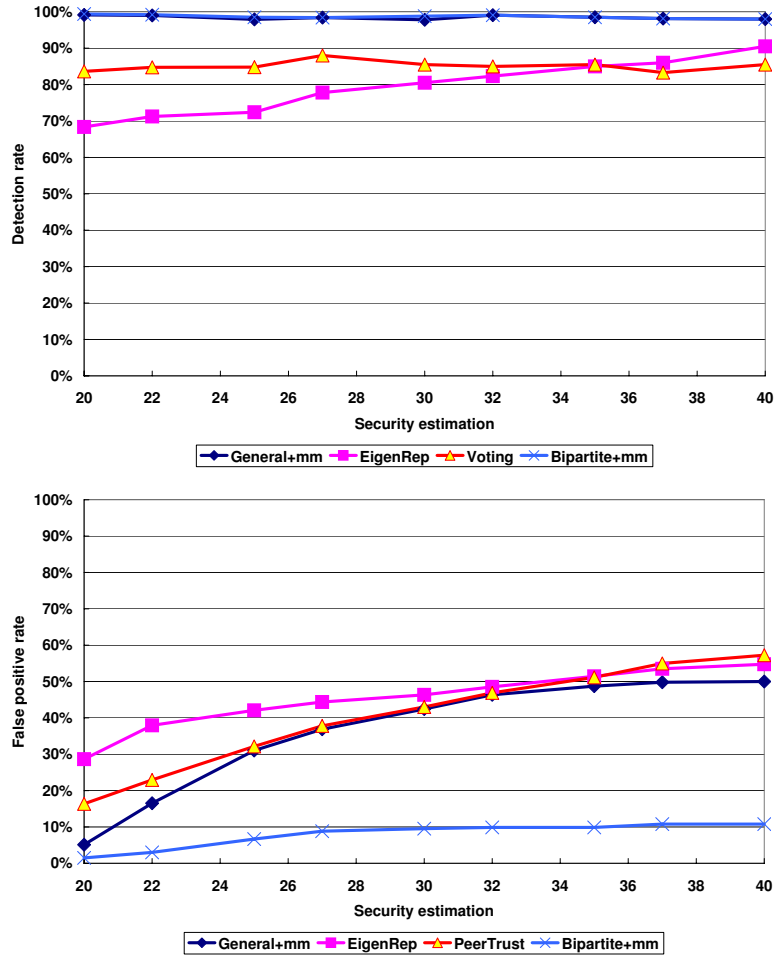


Fig. 7. The impact of the accuracy of security estimation

positives. When the accuracy of the security estimation decreases, the effectiveness of the first phase also decreases. More compromised nodes are instead identified in the second phase by the maximum matching approach. That explains why the false positives increase while detection rates remain high. The detection rates of EigenTrust and PeerTrust in fact improve a little bit when the security estimation accuracy decreases. This is because they always identify the K nodes with the lowest trust values as compromised, which will include more compromised nodes as K increases. But this also increases their false positive rates. For bipartite+mm, even when K is set as twice of the real number of compromised nodes, the false positive rate is still not high. This is because when we reason on the set of alerts and build the inferred graph, we take into consideration not only the nodes which are directly involved in the alerts, but also some other nodes which do not involve in any alerts but can observe those involved nodes. When having more nodes, chances are there will be more nodes in each partition of the bipartite graph, thus the judgement condition $|a_i| + \sum_{j \neq i} |b_j| > K$ in algorithm CollusionCompromisedCore is more easily to be satisfied than the condition $n_s + m > K$ in algorithm AppCompromisedCore.

In summary, our experiments show that the proposed two-phase approach in most cases achieve high detection rates with low false positives, even when sensor nodes are relatively loosely deployed, and compromised nodes form strong local majorities. Also, since they are designed to accommodate the uniqueness of sensor networks, they consistently outperform EigenTrust and PeerTrust, two well-known reputation-based trust management schemes for general decentralized systems, as well as the simple voting mechanism. This is especially true when compromised nodes form strong local majorities.

Our experiments also indicates how alerts can be implicitly inferred from application data for some specific applications as in our example, thus does not introduce any additional communication or computation cost at the network.

In the next experiment, we evaluate the impact on the effectiveness of our algorithms when the security estimation is in fact less than the actual number of compromised nodes in a network. Same as the previous experiment, we set the number of compromised nodes to be 20 and vary the security estimation K from 1 to 19. We observe that our algorithm always throws an exception, stating the size of the minimum vertex cover of the inferred graph is greater than K . This exception indicates that we underestimate the number of compromised nodes when setting K because it contradicts with lemma 3.6. We realize that such underestimation may not always be detected. It is possible that with certain observability graphs, an attacker might be able to compromise more than K sensor nodes, and submit false data and alerts in a way such that the size of the minimum vertex cover of the resulting inferred graph is no more than K . We will study the properties of such graphs in future work. Currently, we treat the exceptions as chances to use multiple K for the detection. That is, if there is an exception, then we should go back and double check our estimation about K , then recompute another bigger K and run our algorithm again.

4.5 Network Evolution

After some compromised nodes are identified and removed, there will be fewer nodes remaining in the system, which will inevitably affect the effectiveness of our scheme. In this section, we conduct experiments to evaluate the degradation of proposed scheme when an attacker continues to compromise more sensor nodes. Specifically, we set the number of sensor nodes in the area to be 200 initially. In each time window, the attacker compromises 15 more nodes. These nodes are chosen from a random new center (x, y) within the area, with concentration of 15. We set the security estimation $K = 20$ to reflect our knowledge that no more than 20 nodes will be compromised within each time window. When the total number of identified compromised nodes is bigger than K , an exception will be thrown to indicate that we need to reevaluate our security estimation. We run both general+mm and bipartite+mm. Figure 8 shows the detection rate and false positive rate at the end of each time window when the base station uses the two algorithms to identify compromised nodes.

We see that as time goes on, the detection rate of both general+mm and bipartite+mm decrease, and the false positive rates increase. This is expected as the density of the network gets smaller when more and more nodes are excluded from the system. To illustrate this, the number of nodes remaining in the network after we exclude those identified as compromised are shown in figure 8(c). As we have shown in section 4.3, when the density of a network decreases, a node will have fewer observers, which make it harder to identify a compromised node definitely.

The result of the experiment suggests the importance of network reconfiguration. Once

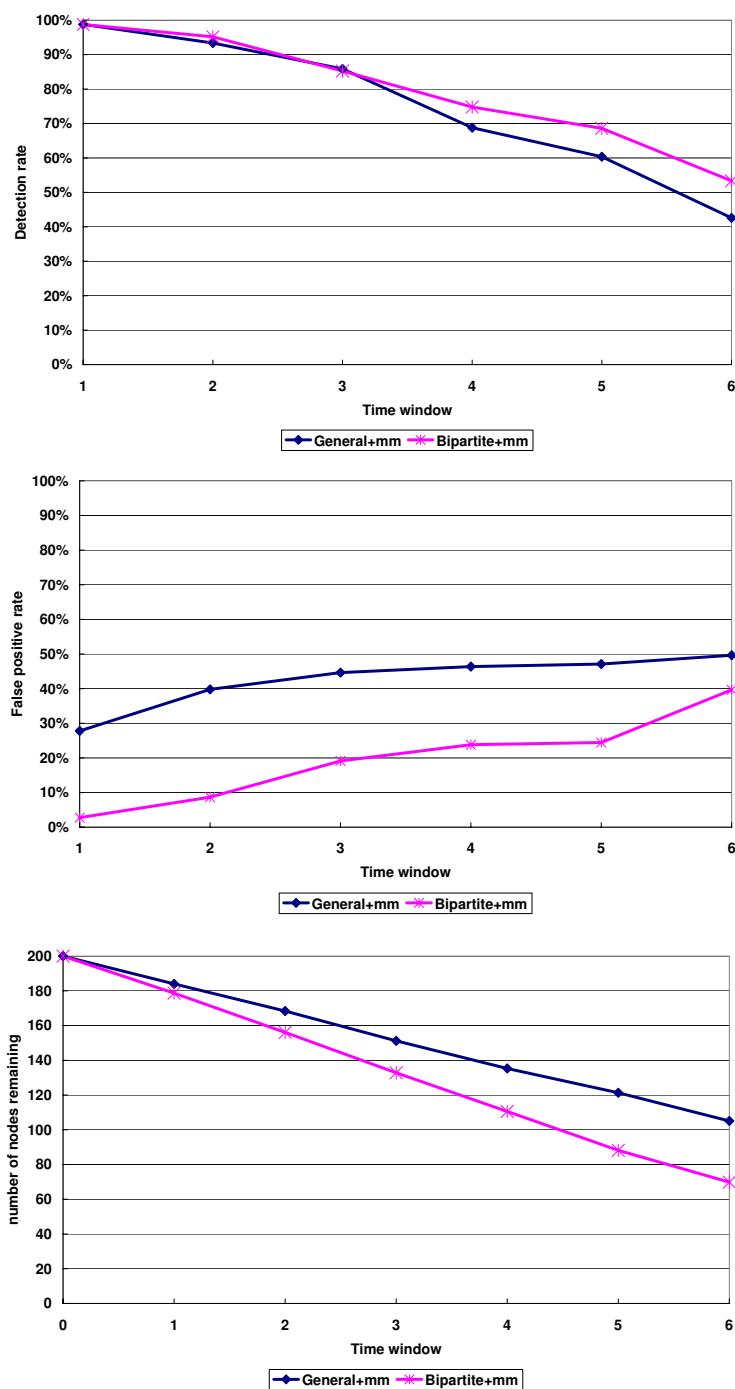


Fig. 8. The effectiveness of the proposed scheme when an attacker continues to compromise more nodes

some compromised nodes are removed, it is necessary to deploy new sensor nodes into the network so that each node is observed by a sufficient number of other nodes. Otherwise, an attacker can selectively compromise those nodes which are insufficiently observed, and avoid being identified. We briefly discuss network reconfiguration in section 5.

5. DISCUSSION

5.1 Attacker's Strategy

The identification algorithm is not secret. Thus it is important to investigate that, with the knowledge of the identification algorithm, what an attacker should do to compromise nodes and maximize its influence in the sensor network. The measure of the influence of compromised nodes often depends on the functionality of specific applications. In this paper, we adopt a simple definition based on the following observation. After some nodes are identified as compromised, they will be removed from the network. The data collected by them will not be used in a sensor network application. Thus, after the identification algorithm is used, the bigger portion of sensor nodes are controlled by an attacker, the larger influence the attacker has.

Specifically, a compromised node may have the following malicious activities: sending false data, sending false alerts against an uncompromised node, and not sending alerts when it should according to a sensor network's detection mechanism. A compromised node may take multiple malicious activities during a period of time. The goal of the attacker is to maximize false information introduced into the systems. In other words, the more nodes sending false data are not detected, the larger influence the attacker has in the sensor network.

If a compromised node s_1 sends false data and one of its observers s_2 is not compromised, then s_2 will issue alerts against s_1 to the base station, and $\{s_1, s_2\}$ forms a suspicious pair in the inferred graph. Due to the maximum matching in the second phase of the identification algorithm, s_1 and s_2 may both be identified as compromised. Thus, in order to be sure that s_1 is definitely not identified as compromised by our identification algorithm, the attacker also needs to compromise all the observers of s_1 . In some sense, the compromised observers of s_1 serve as the protectors of s_1 so that its false data can be accepted by the base station. The protectors may still send legitimate data to the base station. They just do not issue alerts against s_1 . By doing so, neither s_1 nor its observers will be included in any suspicious pairs. Therefore, s_1 will definitely be treated as uncompromised by the identification algorithm.

Generally, we have the following definition.

Definition 5.1. Let G be an observability graph and S be a set of vertices in G . We say $S' \subseteq S$ is a *protected set* of S if for any $s \in S'$, all the observers of s belong to S . The *protection capability* of S is the size of the maximum protected set of S .

Intuitively, the protection capability of S corresponds to the maximum amount of false data that an attacker can definitely inject to the base station once S are compromised.

Given an observability graph G , the attacker may face the following two types of problems. First, suppose the attacker only has the resources to compromise n nodes. Then how to determine a set of n nodes in G so that they have the maximum protection capability. Second, suppose the attacker would like to inject false data from at least t nodes without being detected, then what is the minimum number of sensor nodes it has to compromise.

THEOREM 5.2. *The problem to maximize an attacker's protection capability after compromising n nodes is NP-complete.*

PROOF. First let us look at the graph cutting problem [Marx 2004]: give a graph $G(V, E)$ and two integers k and l , is there a partition $V = X \cup S \cup Y$ such that $|X| = l$, $|S| \leq k$ and there is no edge between X and Y ? This problem is equivalent to the graph cutting problem where $l = |V| - n$. It has been shown by Marx that graph cutting problem is NP-complete by reducing the minimum vertex covering problem to this problem [Marx 2004]. \square

THEOREM 5.3. *The problem to minimize the number of compromised nodes to achieve protection capability t is NP-complete.*

PROOF. Similar to the proof above, this problem is equivalent to the graph cutting problem where $l = t$. \square

The above results are bad news for attacker, fortunately. With limited resources, it is very hard for attackers to inject maximum false information. Also, in order to inject certain amount of false information, the attacker may have to end up compromising more nodes than necessary, which requires more efforts. On the other hand, it is not clear yet whether there exists good approximation algorithm to the above problems. This is one of the topics that we will investigate in our future work.

5.2 Attacks

The proposed framework and algorithms focus on identifying compromised nodes through reasoning about alerts between sensor nodes. One possible attack is that an attacker may repetitively trigger events that can only be monitored by a node s . Thus, the data reported by s may be significantly different from that of others. This may cause many alerts against s , and have s identified as compromised. The essential reason for such attack is that the detection mechanism cannot tell the difference of information of a real phenomenon from pure bogus information generated by a compromised node. This problem cannot be handled by the general framework. Instead, it requires more accurate application detection mechanisms, e.g., having sensor nodes more densely deployed so that any event can be monitored by multiple nodes.

5.3 Network Reconfiguration

Once compromised nodes are identified, it calls for mechanisms to effectively mitigate their impacts. One straightforward method is to remove those sensor nodes from the network through, e.g., key revocation. In some situation, however, this may not be enough. Sensors may provide several services in a single application, such as routing, data sensing and data aggregation. Removing sensors from a network also removes services from them, which may have a significant impact on the functionality of an application. Moreover, as shown by the experiment in section 4.5, when there are fewer nodes in the area, it is also more vulnerable to attacks. Therefore, it is often necessary to reconfigure the service-providing relation between existing sensors, or to deploy new sensors. Sensor network reconfiguration has several goals, including preserving the functionality of a sensor network, minimizing reconfiguration costs and improving the overall trustworthiness of a sensor network. A suitable reconfiguration cost model for sensor networks is essential to achieve the above goals.

5.4 Decentralized approaches

In this paper, we adopt a centralized approach, i.e., the base station collects alerts and identifies potentially compromised nodes. A centralized approach usually offers better accuracy in identifying compromised and malfunctional nodes, since it has a global view of the network. A decentralized approach (as in [Ganeriwal and Srivastava 2004]) is a possible alternative, which limits alerts to be exchanged between nearby nodes. A decentralized approach may incur less communication costs. But without global information, it is in general more difficult to deal with the local majority and collusion of compromised nodes. How to design light-weighted decentralized approach to accurately identify compromised nodes instead of just tolerating them is a challenging problem.

6. RELATED WORK

Much work has been done to provide security primitives for wireless sensor networks, including practical key management [Chan et al. 2003; Du et al. 2003b; Eschenauer and Gligor 2002; Liu and Ning 2003], broadcast authentication [Liu et al. 2003; Perrig et al. 2000; Perrig et al. 2001], and data authentication [Hu and Evans 2003; Przydatek et al. 2003; Zhu et al. 2004] as well as secure in-network processing [Deng et al. 2003]. The work of this paper is complementary to the above techniques, and can be combined to achieve high information assurance for sensor network applications. Several approaches have been proposed to detect and tolerate false information from compromised sensor nodes [Du et al. 2003a; Hu and Evans 2003; Przydatek et al. 2003] through e.g., sampling and redundancy. But they do not provide mechanisms to accurately identify compromised sensor nodes, which is the focus of this paper.

Reputation-based trust management has been studied in different application contexts, including P2P systems, wireless ad hoc networks, social networks and the Semantic Web [Aberer and Despotovic 2001; Kamvar et al. 2003; Lee et al. 2003; Mui et al. 2002; Richardson et al. 2003]. Many trust inference schemes have been proposed. They differ greatly in inference methodologies, complexity and accuracy. As discussed early, the interaction model and assumptions in the above applications are different from sensor networks. Directly applying existing trust inference schemes may not yield satisfactory results in sensor networks.

Ganeriwal et al. [Ganeriwal and Srivastava 2004], propose to detect abnormal routers in sensor networks through reputation mechanism. It is assumed that a sensor's routing quality can be observed by nearby sensors through a watchdog mechanism. Ganeriwal et al. adopt a decentralized trust inference approach. Sensors evaluate each other's trustworthiness by acquiring feedback information from nearby sensors. Ganeriwal et al.'s work shows the usefulness of reputation in sensor networks. But their approach treats a sensor network the same as a typical P2P system, and thus does not capture the unique properties of sensor networks. Further, their work focuses on avoiding services from potentially compromised sensors instead of identifying and excluding them from sensor networks. Further, their work is application specific, and cannot be easily applied to other sensor network applications.

The problem of detecting faulty nodes in multi-processor systems has been studied for a long time. The pioneering work is the PMC model proposed in [Preparata et al. 1967]. Efficient diagnosing algorithms to identify faulty nodes are also proposed [Araki and Shibata 2003; Sullivan 1988; Dahbura and Masson 1984; Fuhrman 1996]. However, all these

algorithms assume that the test assignments follow some special topologies to ensure the system is t -diagnosable in the first place, which means all faulty nodes in the system can be identified as long as there are at most t faulty nodes within it. These algorithms cannot be applied to sensor networks due to their topology requirements, as deployments in sensor networks are often ad hoc.

Some variants of the PMC relaxed from permanent faults to intermittent faults [Dahbura et al. 1987; Kozłowski and Krawczyk 1991], but they need certain assumptions, i.e., the faulty nodes will be faulty following certain probabilities [Dahbura et al. 1987], or the number of incorrect outcomes is bounded [Kozłowski and Krawczyk 1991]. In the most general cases, it is shown by Dahbura and Masson that the problem is NP-complete [Dahbura and Masson 1983b; 1983a].

Previous work on Byzantine fault detection generally focuses on the designing of communication protocols or message/detector structures, so that a proper voting mechanism can lead to the exposure of Byzantine generals [Lamport et al. 1982; Ho et al. 2004]. Sensor nodes are often randomly deployed, thus solutions in this area are not applicable in sensor networks either.

7. CONCLUSION

Wireless sensor networks are often deployed in open environments in an unattended manner. Sensors are subject to capture by attackers, and thus more likely to be compromised. Once keying materials are recovered, an attacker may be able to impersonate compromised nodes completely, and inject false information into sensor network to influence the outcome of an application.

In this paper, we present a general framework that abstracts the essential properties of sensor networks for the identification of compromised sensor nodes. The framework is application-independent, and thus can model a large range of sensor network applications. Built on the alert-based detection mechanisms provided by applications, our framework does not introduce additional communication and computation costs to the network. Based on the framework, we develop efficient algorithms that achieve maximum accuracy without introducing false positives. We further propose techniques to trade off accuracy for increasing the identification of compromised nodes. The effectiveness of these techniques are shown through theoretical analysis and detailed experiments. To the best of our knowledge, our work is the first in the field to provide an application-independent approach to identify compromised nodes in sensor networks. Our algorithm maintains good performances even if we do not have a good estimation of the secrecy of the system.

We plan to extend this work in the following directions. First, we are interested in designing a cost model for sensor network reconfiguration to mitigate the effect of compromised nodes. The model should include possible reconfiguration mechanisms, and consider the multiple functionalities provided by a sensor network and their dependency. Second, we plan to investigate light-weight decentralized approaches, and systematically analyze its benefits and inherent weakness when compared with centralized approaches. Third, we also plan to explore the Bayesian model reasoning, and assign a probability for each edge to infer its likelihood to be abnormal, rather than the current binary model.

REFERENCES

- ABERER, K. AND DESPOTOVIC, Z. 2001. Managing Trust in a Peer-2-Peer Information System. In *Proceedings of the Ninth International Conference on Information and Knowledge Management (CIKM)*.

- ARAKI, T. AND SHIBATA, Y. 2003. (t, k) -diagnosable system: A generalization of the pmc models. *IEEE Trans. on Computers* 52, 7.
- BOSE, P., MORIN, P., STOJIMENOVIC, I., AND URRUTIA, J. 2001. Routing with guaranteed delivery in ad hoc wireless networks. *ACM Wireless Networks* 7, 6, 609–616.
- CAMTEPE, S. AND YENER, B. 2004. Combinatorial design of key distribution mechanisms for wireless sensor networks. In *9th European Symposium On Research in Computer Security (ESORICS'04)*.
- CHAN, H., PERRIG, A., AND SONG, D. 2003. Random key predistribution schemes for sensor networks. In *IEEE Symposium on Security and Privacy (SP'03)*.
- CROSSBOW TECHNOLOGY INC. 2003. *MTS/MDA Sensor and Data Acquisition Boards User Manual*.
- DAHURA, A. AND MASSON, G. 1983a. Greedy diagnosis of an intermittent-fault/transient-upset tolerant system design. *IEEE Trans. on Computers C-32*, 10, 953–957.
- DAHURA, A. AND MASSON, G. 1983b. Greedy diagnosis of hybrid fault situations. *IEEE Trans. on Computers C-32*, 8, 777–782.
- DAHURA, A. AND MASSON, G. 1984. An $o(n^{2.5})$ fault identification algorithm for diagnosable systems. *IEEE Trans. on Computers C-33*, 6, 486–492.
- DAHURA, A., SABNANI, K., AND KING, L. 1987. The comparison approach to multiprocessor fault diagnosis. *IEEE Trans. on Computers C-36*, 3, 373–378.
- DENG, J., HAN, R., AND MISHRA, S. 2003. Security support for in-network processing in wireless sensor networks. In *2003 ACM Workshop on Security in Ad Hoc and Sensor Networks (SASN '03)*.
- DENG, J., HAN, R., AND MISHRA, S. 2004. A Robust and Light-Weight Routing Mechanism for Wireless Sensor Networks. In *Workshop on Dependability Issues in Wireless Ad Hoc Networks and Sensor Networks (DIWANS)*.
- DU, W., DENG, J., HAN, Y. S., AND VARSHNEY, P. K. 2003a. A Witness-Based Approach For Data Fusion Assurance In Wireless Sensor Networks. In *IEEE 2003 Global Communications Conference (GLOBECOM)*.
- DU, W., DENG, J., HAN, Y. S., AND VARSHNEY, P. K. 2003b. A pairwise key pre-distribution scheme for wireless sensor networks. In *10th ACM Conference on Computer and Communications Security (CCS'03)*.
- DU, W., FANG, L., AND NING, P. 2005. Lad: Localization anomaly detection for wireless sensor networks. In *19th IEEE International Parallel and Distributed Processing Symposium (IPDPS'05)*.
- ESCHENAUER, L. AND GLIGOR, V. D. 2002. A key-management scheme for distributed sensor networks. In *9th ACM conference on Computer and communications security (CCS '02)*.
- FUHRMAN, C. P. 1996. Comparison-based diagnosis in faulttolerant, multiprocessor systems. Ph.D. thesis, Department of Computer Science, Swiss Federal Institute of Technology in Lausanne (EPFL).
- GANERIWAL, S. AND SRIVASTAVA, M. B. 2004. Reputation-based Framework for High Integrity Sensor Networks . In *ACM Security for Ad-hoc and Sensor Networks (SASN 2004)*.
- GOLBECK, J. AND HENDLER, J. 2004. Accuracy of Metrics for Inferring Trust and Reputation in Semantic Web-based Social Networks. In *International Conference on Knowledge Engineering and knowledge Management (EKAW)*. Northamptonshire, UK.
- HO, T., LEONG, B., KOETTER, R., MEDARD, M., EFFROS, M., AND KARGER, D. 2004. Byzantine Modification Detection in Multicast Networks using Randomized Network Coding. In *2004 IEEE International Symposium on Information Theory (ISIT)*.
- HU, L. AND EVANS, D. 2003. Secure aggregation for wireless networks. In *Workshop on Security and Assurance in Ad Hoc Networks*.
- KAMVAR, S., SCHLOSSER, M., AND GARCIA-MOLINA, H. 2003. EigenRep: Reputation Management in P2P Networks. In *Twelfth International World Wide Web Conference*.
- KOZŁOWSKI, W. AND KRAWCZYK, H. 1991. A comparison-based approach to multicomputer system diagnosis in hybrid fault situations. *IEEE Trans. on Computers C-40*, 11, 1283–1287.
- LAMPORT, L., SHOSTAK, R., AND PEASE, M. 1982. The byzantine generals problem. *ACM Transactions on Programming Languages and Systems* 4, 3.
- LAWRENCE, R., SERGEY, B., RAJEEV, M., AND TERRY, W. 1998. The PageRank Citation Ranking: Bringing Order to the Web. Tech. rep., Department of Computer Science, Stanford University.
- LEE, S., SHERWOOD, R., AND BHATTACHARJEE, B. 2003. Cooperative Peer Groups in NICE. In *INFOCOM*.
- LIU, D. AND NING, P. 2003. Establishing pairwise keys in distributed sensor networks. In *10th ACM conference on computer and communications security (CCS '03)*.

- LIU, D., NING, P., AND DU, W. 2003. Efficient distribution of key chain commitments for broadcast authentication in distributed sensor networks. In *10th Annual Network and Distributed System Security Symposium (NDSS'03)*.
- LIU, D., NING, P., AND DU, W. 2005. Detecting Malicious Beacon Nodes for Secure Location Discovery in Wireless Sensor Networks. In *Proceedings of the The 25th International Conference on Distributed Computing Systems (ICDCS '05)*.
- LIU, D., NING, P., AND LI, R. 2004. Establishing pairwise keys in distributed sensor networks. *ACM Transactions on Information and System Security*.
- MARX, D. 2004. Parameterized complexity of constraint satisfaction problems. In *19th Annual IEEE conference on computational complexity*.
- MICALI, S. AND VAZIRANI, V. 1980. An $o(\sqrt{|V|})|E|$ algorithm for finding maximum matchings in general graphs. In *21st. Symp. Foundations of Computing*.
- MUI, L., MOHTASHEMI, M., AND HALBERSTADT, A. 2002. A Computational Model of Trust and Reputation. In *35th Hawaii International Conference on System Science*.
- PERRIG, A., CANETTI, R., SONG, D., AND TYGAR, D. 2000. Efficient authentication and signing of multicast streams over lossy channels. In *IEEE Symposium on Security and Privacy*.
- PERRIG, A., SZEWCZYK, R., WEN, V., CULLER, D., AND TYGAR, J. D. 2001. SPINS: Security Protocols for Sensor Networks. In *Seventh Annual ACM International Conference on Mobile Computing and Networks (MobiCom 2001)*.
- PREPARATA, F. P., METZE, G., AND CHIEN, R. T. 1967. On the connection assignment problem of diagnosable systems. *IEEE Trans. on Electronic Computers* 16, 6, 848–854.
- PRZYDATEK, B., SONG, D., AND PERRIG, A. 2003. SIA: Secure information aggregation in sensor networks. In *First ACM Conference on Embedded Networked Sensor Systems (SenSys'03)*.
- RICHARDSON, M., AGRAWAL, R., AND DOMINGOS, P. 2003. Trust Management for the Semantic Web. In *Proceedings of the Second International Semantic Web Conference*.
- SULLIVAN, G. F. 1988. An $o(t^3 + |E|)$ fault identification algorithm for diagnosable systems. *IEEE Trans. on Computers* 37, 4.
- VAZIRANI, V. V., Ed. 2001. *Approximation Algorithms*. Springer-Verlag.
- XIONG, L. AND LIU, L. 2002. Building Trust in Decentralized Peer-to-Peer Electronic Communities. In *The 5th International Conference on Electronic Commerce Research (ICECR)*.
- YE, F., LUO, H., LU, S., AND ZHANG, L. 2004. Statistical en-route filtering of injected false data in sensor networks. In *IEEE INFOCOM*.
- YU, B. AND SINGH, M. P. 2002. An Evidential Model of Distributed Reputation Management. In *Proceedings of the 1st International Joint Conference on Autonomous Agents and MultiAgent Systems (AAMAS)*.
- ZHANG, Q., YU, T., AND NING, P. 2006. A Framework for Identifying Compromised Nodes in Sensor Networks. In *2nd IEEE Communications Society/CreateNet International Conference on Security and Privacy in Communication Networks (SecureComm 2006)*.
- ZHU, S., SETIA, S., JAJODIA, S., AND NING, P. 2004. An Interleaved Hop-by-Hop Authentication Scheme for Filtering of Injected False Data in Sensor Networks. In *IEEE Symposium on Security and Privacy*, pages 260–272.

Appendix

Here we provide the solution of the distribution of d in the sensor behavior model of in section 2.4. The problem can be interested as: Let X and Y are two independent random variables with standard normal distribution $N(0, \sigma^2)$, Let $U = \sqrt{X^2 + Y^2}$, what is the distribution for U ?

Solution: This is a bivariate transformation problem, and we have following fact:

Let (X, Y) be a bivariate *continuous* random vector with a known probability distribution. Now consider a new bivariate random vector (U, V) defined by $U = g_1(X, Y)$ and $V = g_2(X, Y)$, where $g_1(x, y)$ and $g_2(x, y)$ are some specified functions. If (X, Y) has

probability density function(pdf) $f_{X,Y}(x, y)$, then the joint pdf of (U, V) can be expressed in terms of $f_{X,Y}(x, y)$. We assume that $u = g_1(x, y)$ and $v = g_2(x, y)$ defines a one-to-one transformation, then we can derive the inverse transformation $x = h_1(u, v)$ and $y = h_2(u, v)$. So the joint density function of random vector (U, V) is

$$f_{U,V}(u, v) = f_{X,Y}(h_1(u, v), h_2(u, v))|J|$$

where

$$f_{X,Y}(x, y) = \frac{1}{2\pi\sigma^2} e^{-x^2/2\sigma^2} e^{-y^2/2\sigma^2}$$

$$J = \begin{vmatrix} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{vmatrix} \quad x = h_1(u, v), y = h_2(u, v)$$

This is just a simple conclusion, you can see some details from any fundamental statistics book.

Now back to our problem, we are only interested in the distribution of $U = \sqrt{X^2 + Y^2}$, but we can construct a function $V = g(X, Y)$, and use the conclusion above to get the joint density function of (U, V) , then find the marginal probability of U , $f_U(u)$, which is what we want.

Suppose our new bivariate random vector is

$$U = g_1(X, Y) = \sqrt{X^2 + Y^2} \quad V = g_2(X, Y) = Y$$

However, g_1 and g_2 are not one-to-one transformation since the point (x, y) and $(-x, y)$ are both mapped into the same (u, v) , we can not use the fact above directly. But if we restrict consideration to either positive or negative values of x , then the transformation is one-to-one. Let

$$A_1 = \{(x, y) : x > 0\} \quad A_2 = \{(x, y) : x < 0\}$$

then the inverse transformation for A_1 is

$$x = \sqrt{u^2 - v^2}, \quad y = v$$

for A_2 is

$$x = -\sqrt{u^2 - v^2}, \quad y = v$$

The Jacobians from the two inverses are $J_1 = J_2 = |u|/\sqrt{u^2 - v^2}$, therefore, the joint density function of (U, V) is

$$\begin{aligned} f_{U,V}(u, v) &= \frac{1}{2\pi\sigma^2} e^{-(u^2-v^2)/2\sigma^2} e^{-v^2/2\sigma^2} \frac{|u|}{\sqrt{u^2 - v^2}} \\ &+ \frac{1}{2\pi\sigma^2} e^{-(u^2-v^2)/2\sigma^2} e^{-v^2/2\sigma^2} \frac{|u|}{\sqrt{u^2 - v^2}} \\ &= \frac{1}{\pi\sigma^2} e^{-u^2/2\sigma^2} \frac{|u|}{\sqrt{u^2 - v^2}}, \quad 0 \leq u < \infty, 0 \leq v \leq u \end{aligned}$$

From this the marginal pdf of U can be computed to be

$$\begin{aligned}
 f_U(u) &= \int_0^u \frac{1}{\pi\sigma^2} e^{-u^2/2\sigma^2} \frac{u}{\sqrt{u^2-v^2}} dv \\
 &= \frac{e^{-u^2/2\sigma^2}}{\pi\sigma^2} \int_0^u \frac{u}{\sqrt{u^2-v^2}} dv \\
 &= \frac{u}{2\pi^2\sigma^2} e^{-u^2/2\sigma^2}
 \end{aligned}$$

where the integral

$$\int_0^u \frac{u}{\sqrt{u^2-v^2}} dv = \frac{u}{2\pi}$$