

# A Hybrid Intrusion Detection System for IoT Applications with Constrained Resources

Chao Wu, Chongqing Vehicle Test & Research Institute Co. Ltd., Chongqing, China

Yuan'an Liu, School of Electronic Engineering, Beijing University of Posts and Telecommunications, Beijing, China

Fan Wu, School of Electronic Engineering, Beijing University of Posts and Telecommunications, Beijing, China

Feng Liu, SKLOIS, IIE & SCS UCAS, CAS, Beijing, China

Hui Lu, Institute of Microelectronics of the Chinese Academy of Sciences, Beijing, China

Wenhao Fan, School of Electronic Engineering, Beijing University of Posts and Telecommunications, Beijing, China

Bihua Tang, School of Electronic Engineering, Beijing University of Posts and Telecommunications, Beijing, China

## ABSTRACT

Network security and network forensics technologies for the Internet of Things (IoT) need special consideration due to resource-constraints. Cybercrimes conducted in IoT focus on network information and energy sources. Graph theory is adopted to analyze the IoT network and a hybrid Intrusion Detection System (IDS) is proposed. The hybrid IDS consists of Centralized and Active Malicious Node Detection (CAMD) and Distributed and Passive EEA (Energy Exhaustion Attack) Resistance (DPER). CAMD is integrated in the genetic algorithm-based data gathering scheme. CAMD detects malicious nodes manipulated by cyber criminals and provides digital evidence for forensics. DPER is implemented in a set of communication protocols to alleviate the impact of EEA attacks. Simulation experiments conducted on NS-3 platform showed the hybrid IDS proposed detected and traced malicious nodes precisely without compromising energy efficiency. Besides, the impact of EEA attacks conducted by cyber criminals was effectively alleviated.

## KEYWORDS

Cybercrime, Energy Efficiency, Genetic Algorithm, Graph Theory, Internet Of Things, Network Forensics

## INTRODUCTION

Network forensics is the reconstruction of network event to provide definitive insight into action and behavior of users, applications as well as devices (Schwartz, 2010). Network forensics technologies focus on recording evidence of a network attack (Adeyemi, Razak, & Nor Azhan, 2013). However, Internet of Things (IoT) is a special network which integrates sensors and other objects to connect everything in our life together. The information in IoT is usually privacy-sensitive and even confidential, so IoT will become the objective of cyber criminals (Alaba, Othman, Hashem, & Alotaibi, 2017). Due to the device miniaturization and energy-efficiency of IoT, traditional network forensics technologies are not suitable for IoT. Thus, the network forensics technologies specialized for cybercrimes aiming at IoT are of great importance and challenging in the era of IoT. Different from traditional computer networks, IoT networks are typically Low-power and Lossy Networks (LLN) (Teklemariam, Van Den Abeele, & et al, 2016), so energy efficiency must be taken into consideration when it comes to network security and network forensics technology designs for IoT.

DOI: 10.4018/IJDCF.2020010106

This article, originally published under IGI Global's copyright on January 1, 2020 will proceed with publication as an Open Access article starting on January 27, 2021 in the gold Open Access journal, International Journal of Digital Crime and Forensics (converted to gold Open Access January 1, 2021), and will be distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>) which permits unrestricted use, distribution, and production in any medium, provided the author of the original work and original publication source are properly credited.

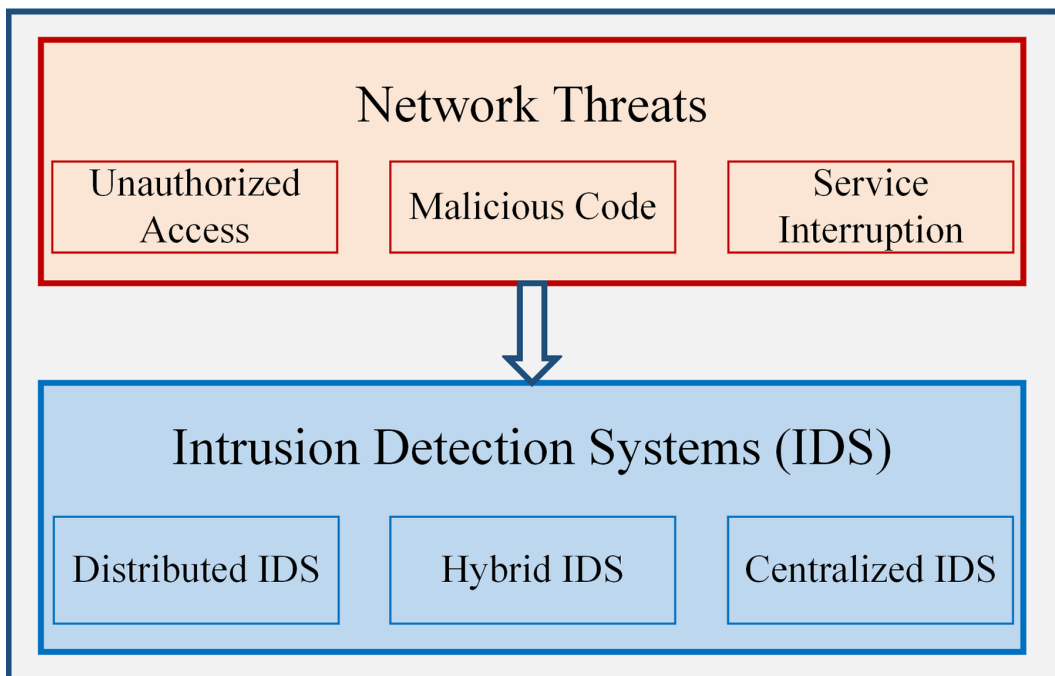
Intrusion Detection Systems (IDSs) can be categorized into three types by placement (Zarpelao, Miani, Kawakani, & de Alvarenga, 2017), as shown in Figure 1. Distributed IDS mean the detection system is placed in every physical node. Distributed IDSs are suitable for smart devices with higher computational capability and energy sources. Correspondingly, centralized IDS only rely on single or several dedicated components in the network to complete the detection work. Hybrid IDS combines distributed and centralized technologies to get the job done.

Aiming at computer networks, threats can be categorized into unauthorized access, malicious code and service interruption (Ahmed, 2017) as showed in Figure 1. In IoT networks, cyber criminals may manipulate data nodes in the network illegally, and generate plenty of fake or harmful information. Besides, unauthorized cyber criminals may access data nodes in IoT networks to perform Denial of Service (DoS) attacks. One form of DoS attacks in IoT is Energy Exhaustion Attack (EEA) (Alrajeh, Khan, Lloret, & Loo, 2014). EEA accelerates the expiration of the network lifetime and is fatal to the performance of IoT.

Sink mobility is recognized as an efficient method to improve the performance of IoT. However, mobility-constrained mobile sinks exist in many IoT applications, such as railway-based (Smeets, Shih, Zuniga, Hagemeyer, & Marrón, 2013) or automobile-based (Huang & Savkin, 2016) information collection applications, mountainous or canal environment monitoring applications, and even the information collection application for Smart Grid.

This paper designs an information and energy-related IDS with hybrid mechanism for IoT applications with a path-constrained mobile sink. The hybrid IDS provides a trace-back mechanism for network forensics and enhances the network safety. The main contributions of this paper are summarized as follows:

Figure 1. Network threats and IDS categories



- Graph theory is utilized to describe the IoT network and to group the data nodes with grids. Data nodes contained in the same grid can be regarded as a whole. Network analysis can be simplified by grids.
- The first part of the hybrid IDS is called Centralized and Active Malicious Node Detection (CAMD). CAMD is integrated in the graph-based data gathering scheme of IoT. By a customized genetic algorithm, network nodes which are manipulated by cyber criminals can be identified according to the changes in the data gathering scheme. The changes are got by matrix comparison and the comparison results can be regarded as digital evidence for network forensics.
- A Distributed and Passive EEA Resistance (DPER) mechanism as the second part of the hybrid IDS is fulfilled by a pair of data delivery protocols. DPER is designed to alleviate the impact of EEA on the network performance of the IoT application.
- Simulation experiments are conducted on the Network Simulator-3 (NS-3) platform. In the experiments, cybercrime behaviors including fake information transmitting and EEA attacks are simulated. The experiment results show that the hybrid IDS precisely identifies the cyber criminals who conduct malicious behaviors at data nodes and provides strong digital evidence. Besides, the impact of EEA attacks on the network lifetime is alleviated. The energy efficiency of the IoT network was improved.

## BACKGROUND

With the development of anti-cybercrime technologies, many advanced IDSs have been proposed in recent years, and they were categorized into three types by placement.

Researches on distributed IDSs focused on the energy consumption caused by the massive installation of the IDSs at network nodes. Lightweight IDSs were effective approaches to reduce the extra energy consumption and fulfilled the intrusion detection tasks. Oh proposed the IDS with matching algorithm of malicious behaviors and packet payloads (Oh, Kim, & Ro, 2014). To achieve the lightweight IDS, the authors reduced the iteration times of the pattern matching algorithm. On the other hand, Lee (Lee, Wen, Chang, Chiang, & Hsieh, 2014) discovered another way to reduce the overhead of IDS. The authors established the lightweight IDS under a low energy-consuming communication protocol. However, traditional distributed IDSs required relatively high computational capability for each network node to detect abnormal behaviors while networking.

On the contrary, the intrusion behaviors were detected by dedicated nodes in centralized IDSs. Wallgren utilized heartbeat protocol to fulfill the global monitoring (Wallgren, Raza, & Voigt, 2013) with centralized IDS method. The overheads of heartbeat packets were considerable. For energy-constrained IoT, the massive energy consumption for such global heartbeat packets was not acceptable. Gunasekaran (Gunasekaran & Periakaruppan, 2017) installed an anti-DoSL (Denial of Sleep) system on the base station (BS) in IoT. Based on genetic algorithm, the authors adopted an encryption algorithm with specialized key packets to identify DoSL attacks from malicious nodes. However, the dedicated nodes in centralized IDSs took a while to detect the intrusion behaviors, thus the purposes of cyber criminals might have already been fulfilled before detection and prevention.

Therefore, energy efficient and real-time IDSs specialized for IoT applications were needed desperately. Hybrid IDSs always partitioned the network or dynamically employed nodes to detect malicious behaviors. Amaral (Amaral, Oliveira, Rodrigues, Han, & Shu, 2014) proposed a competition approach that selected robust nodes to monitor adjacent nodes. However, the approach required the selected nodes to be resourceful and the selection consumed extra energy. Le (Le, Loo, Chai, & Aiash, 2016) and Joby (Joby & Sengottuvelan, 2015) utilized clustering method to detect intrusions. Cluster heads (CHs) detected malicious behaviors conducted by member nodes in the cluster. However, cluster heads would deplete rapidly due to extra computation. Alrajeh (Alrajeh, Khan, Lloret, & Loo, 2014) introduced Artificial Neural Network (ANN) into attack detection in IoT. The authors tried to offset the energy consumption caused by the ANN-based IDS and energy harvest method was their choice.

Aiming at the IoT applications with a path constrained mobile sink, this paper proposed a hybrid IDS to trace intrusion sources and to provide digital evidence for forensics. Considering the energy constraint of IoT, the proposed IDS consumed no extra energy.

## SYSTEM MODELING AND PROBLEM ANALYSIS

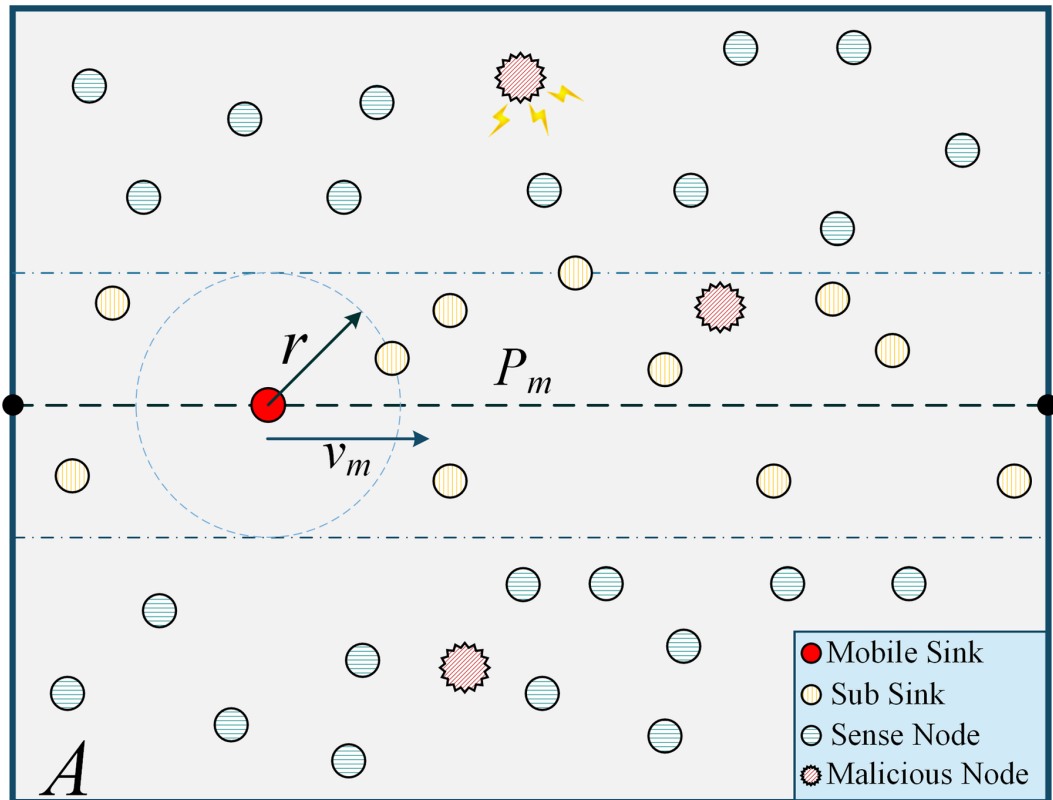
### System Model

Network area  $A$  is a two-dimensional square and  $n$  static data nodes denoted by  $\mathbb{N}$  with initial energy source  $E_{mit}$  are deployed into  $A$  randomly. The movement path of the mobile sink is exactly the symmetric line of  $A$ , denoted by  $P_m$ . So only the upper half of  $A$  is discussed in the remaining content of this paper. Including the mobile sink, all data nodes' max wireless communication radius is  $r$  and the memory of each node is enough for data buffering and routing information recording. The data nodes can locate themselves and record the location information in their memory when deployed. The energy consumption of location can be ignored.

Malicious nodes refer to the data nodes manipulated by cyber criminals while networking. Malicious nodes are with the same physical attributes as normal data nodes. However, due to the simplification of IoT nodes, fake information reporting and EEA attacks are considered as the main malicious behaviors. The system model is illustrated as Figure 2.

With unlimited energy, the mobile sink travels along  $P_m$  with a constant speed  $v_m$  and owns extra computing capability. When the mobile sink gets to the either end point of  $P_m$ , it will turn

Figure 2. The system model



around and head for another end point with the same speed. The data nodes along the mobile sink path that are within the one-hop communication range of the mobile sink are regarded as sub sinks.

Data nodes start to transmit data at a fixed time interval  $t_{int}$  with the data generation rate  $p$  after the destination sub sink is matched. Sub sinks don't transmit data at once but store the data in their buffer along with the data gathered from remote data nodes (Table 1).

### Application of Graph Theory

Graph theory is used in many network issues (Zhou, Du, Shu, Hancke, Niu, & Ning, 2016; Nake & Chatur, 2016; Liu, Zhang, Shen, Fu, & Linge, 2016) and simplifies the analysis of the network. Network area  $A$  is divided into small equal-sized square grids denoted by  $g_i$  as shown in Figure 3.

The edge length of each grid is  $\frac{\sqrt{2}}{2}r$ , so that the data nodes inside a certain grid can communicate with other nodes in the same grid.

After the division, the intrusion detection and energy efficiency problems can be simplified by regarding the nodes in a certain grid as a whole. The grids that contain sub sinks are denoted by  $g_i^{ss}$ . Data nodes contained in the same grid match a sub sink grid  $g_i^{ss}$  and transmit their data to  $g_i^{ss}$ . When data packets arrive at  $g_i^{ss}$ , the sub sinks inside  $g_i^{ss}$  receive the packets evenly according to the

**Table 1. Important symbols and descriptions for the system model**

Symbol	Description
$R_g / C_g$	The number of rows/columns of the matrix of grids in the graph.
$n_{ssg}$	The total amount of sub sink grids. $n_{ssg} = C_g$ .
$g_i$	The symbol that represents the square grid.
$g_i^{ss}$	The symbol represents the grid that contains sub sinks, and $i \in [1, n_{ssg}]$ .
$n_i^{ss}$	The amount of sub sink nodes in sub sink grid $g_i^{ss}$ , $n_i^{ss} =  g_i^{ss} $ .
$n_i^{sng}$	The amount of data node grids assigned to sub sink grid $g_i^{ss}$ .
$n_i^{sn}$	The accumulated amount of data nodes in the data node grids assigned to $g_i^{ss}$ .
$n_{ss}$	Total amount of sub sink nodes. $n_{ss} = \sum_{i=1}^{n_{ssg}} n_i^{ss}$ .
$n_i$	The number of nodes inside some grid $g_i$ , $ g_i  = n_i$ .
$n_g$	The total amount of grids of $G$ , where $ G  = n_g$ .

transmission protocol. Besides, the nodes which can communicate with the nodes in adjacent grids are regarded as pipe nodes.

After the division,  $A$  can be described with a graph, denoted by  $G$ . Mathematically, a graph  $G$  can be represented by its vertices  $V$  and corresponding edges  $E$ , like  $G = (V, E)$  (Zhang, Zhou, Zhao, et al, 2017; Yun, Xia, Behdani, & Smith, 2013; Hasan, AI-Rizzo, & Gunay, 2017), where  $|V| = n_g$ , is the amount of grids and  $|E|$  represents the total amount of pipe nodes in all grids.

The graph operation not only simplifies the network description but also prevents nodes at close locations from unreasonably choosing different sub sinks. Furthermore, benefitted by the division, the computational complexity of the network dropped from  $O(n)$  to  $O(n_g)$ .

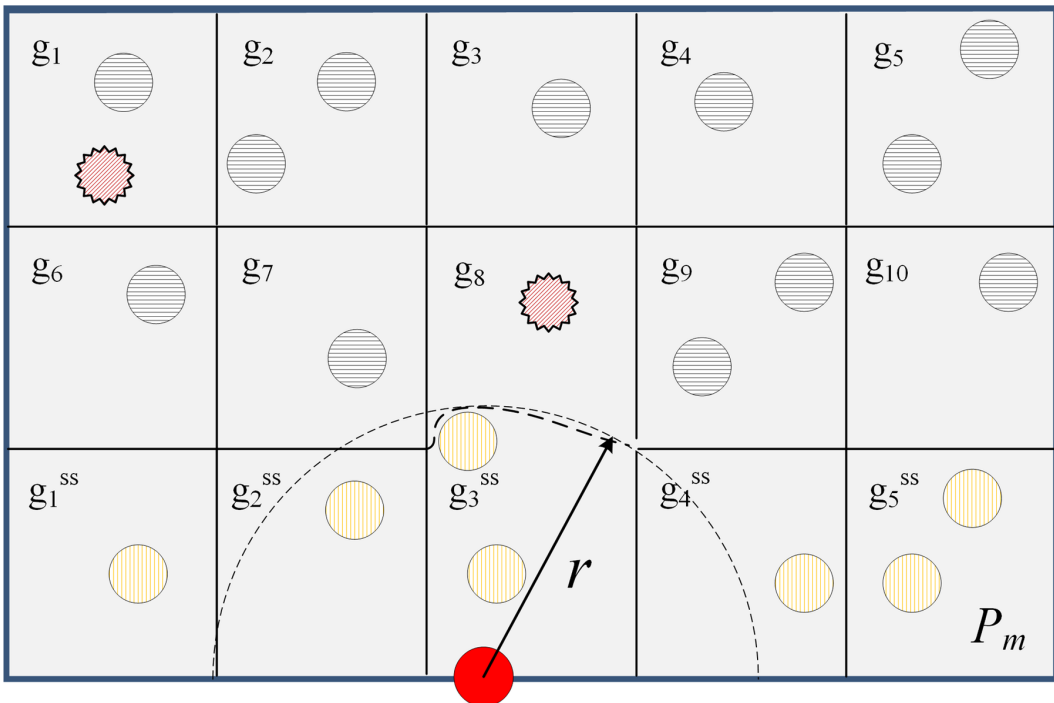
### Analysis of Energy Model

The purpose of the hybrid IDS is to achieve intrusion detection and crime source trace-back without extra energy consumption. Qin (Qin, Hempstead, & Yang, 2009) indicated that, information processing, wireless communication and data aggregation consume most of the energy installed in data nodes.

The transmitting and receiving power can be assumed as fixed for short distance communications, so the power consumption is independent of the transmission distance. Besides, the total energy consumed by data aggregation at each node can be regarded as constants during  $t_{int}$ . The total energy consumed by a data node during  $t_{int}$  can be formulized by:

$$E_{node} \approx e(p_{tx} + p_{rx}) \tag{1}$$

Figure 3. An example of the graphing operation



where  $e$  is a factor indicating the energy consumption per bit for the transmitting and receiving circuits.  $p_{tx}$  and  $p_{rx}$ , respectively, represents the total bytes transmitted and received by a node during  $t_{int}$ .

Furthermore, the total bytes  $p_{tx}^i$  transmitted by node  $i$  can be accounted by:

$$p_{tx}^i = p_{rx}^i + p t_{int} \quad (2)$$

where  $p_{rx}^i$  is the total bytes received by node  $i$ .

Based on Equation (1) and Equation (2), the total energy consumption  $E_{total}$  of the whole network is:

$$E_{total} \approx ept \sum_{i=1}^{n_g} \sum_{k=1}^{n_i} (2h_{ik} + 1)_{int} \quad (3)$$

where  $h_{ik}$  represents the shortest hop from the  $k$  th node in grid  $g_i$  to its matched sub sink grid. The shortest hop is obtained by the routing establishment algorithm discussed later.

As discussed above, any data node grid  $g_i$  chooses but only one reachable sub sink grid as the packet reporting destination of the data nodes inside it. So the total energy consumption  $E_i^{ss}$ , total energy for reception  $E_i^{rx}$  and total energy for transmission  $E_i^{tx}$  in a data interval  $t_{int}$  of sub sinks in a certain sub sink grid  $g_i^{ss}$  are get:

$$E_i^{ss} = E_i^{rx} + E_i^{tx} = ept \left( 2 \sum_{k=1}^{n_i^{msg}} n_{ik} + n_i^{ss} \right) [1, n_{ssg}]_{int} \quad (4)$$

where  $n_i^{ss}$  represents the amount of sub sink nodes in  $g_i^{ss}$ , and  $n_{ik}$  represents the node amount of the  $k$  th data node grid that is matched with  $g_i^{ss}$ .

To measure the equilibrium of energy consumption among all sub sinks, the variance  $D(E_{SS})$  is formulized:

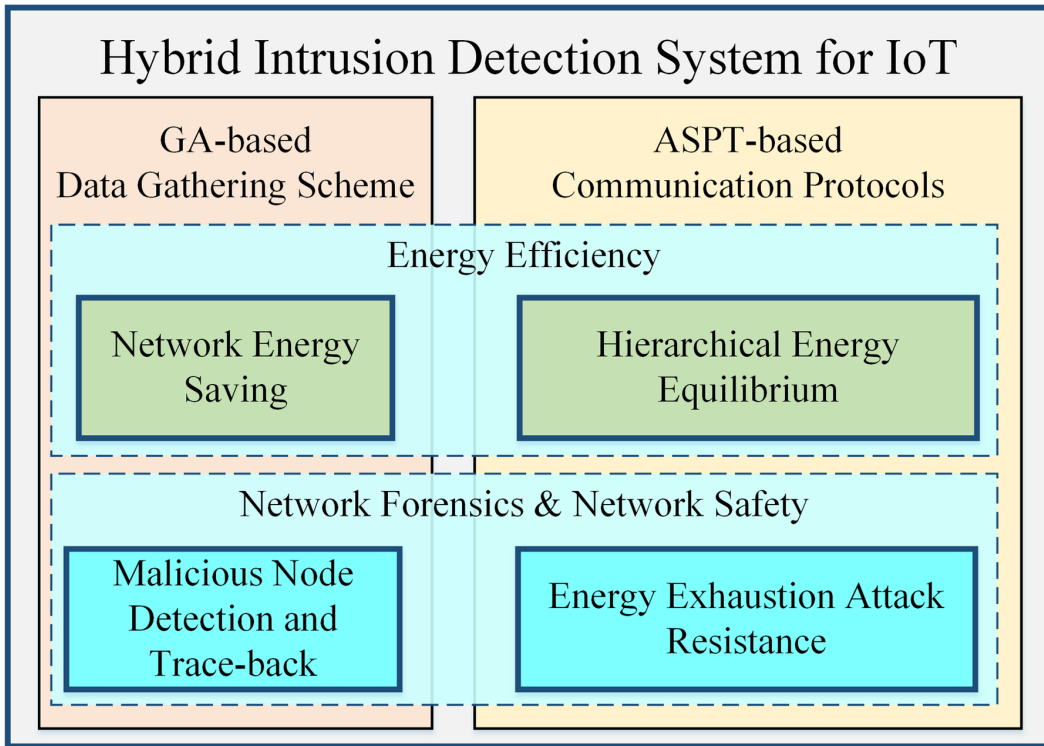
$$D(E_{SS}) \approx \frac{1}{n_{ss}} \sum_{i=1}^{n_{ssg}} \left( \frac{E_i^{ss}}{n_i^{ss}} - \overline{E_{SS}} \right)^2 \cdot n_i^{ss} \quad (5)$$

where  $\overline{E_{SS}} = \frac{1}{n_{ss}} \sum_{i=1}^{n_{ssg}} E_i^{ss}$  is the arithmetic mean value of the energy consumption of all sub sinks in the network.

## Hybrid Intrusion Detection System

Most researches on IDSs designed for traditional computer networks compromised the energy efficiency. However, IoT nodes are not capable for traditional distributed or centralized methods in IDSs due to their simplicity and limited resources. This paper focuses on cybercrime detection and

Figure 4. Illustration of the hybrid Intrusion Detection System proposed



forensics for IoT applications without compromising energy efficiency. A hybrid IDS is integrated in the data gathering function of the IoT network. The hybrid IDS is illustrated in Figure 4.

The hybrid IDS is made up of two primary modules. Genetic algorithm-based (GA) data gathering scheme is proposed to improve the energy efficiency of the network. Malicious behaviors are detected during the data gathering process through a matrix comparison method. The second part of the hybrid IDS works in the form of a set of communication protocols. The protocols are named after Advanced Shortest Path Tree (ASPT). ASPT not only dynamically balances the node energy consumption hierarchically but also alleviates the impact of EEA conducted by the malicious nodes.

## CENTRALIZED AND ACTIVE MALICIOUS NODE DETECTION (CAMD)

### GA-based Data Gathering Scheme

Based on the discussion above, to achieve energy efficiency of the whole network, the value of total energy consumption  $E_{total}$  during each data generation interval  $t_{int}$  needs to be minimized:

$$\min ept \sum_{i=1}^{n_g} \sum_{k=1}^{n_i} (2h_{ik} + 1)_{int} \tag{6}$$

subject to:



$$\min D(E_{SS}) \quad (7)$$

Constraint Equation (7) guarantees the energy consumption among sub sinks is even as far as possible. Considering the idealistic situation that, during each data interval  $t_{int}$  all data packets are transmitted to the destination sub sink grids successfully, which means the total bytes received by all sub sinks during  $t_{int}$  can be regarded as a constant. Therefore, the sum of energy consumed by all sub sinks is accordingly a constant. Then constraint Equation (7) can be transformed into:

$$\min \frac{1}{n_{ssg}} \sum_{i=1}^{n_{ssg}} \left( n_i^{sn} - \overline{n^{sn}} \right)^2 \quad (8)$$

where  $n_i^{sn}$  is the total amount of data nodes in the data grids that are matched with sub sink grid  $g_i^{ss}$  and  $\overline{n^{sn}}$  is accordingly the arithmetic mean value of all  $n_i^{sn}$ .

Constraint Equation (8) means that, to achieve the equilibrium of energy consumption among all sub sinks is approximately equivalent to minimize the variance of the total amounts of data nodes matched with each sub sink grid.

According to the above analysis, the goal of the energy-concerned part of CAMD is to look for an optimal assignment of data grids to sub sink grids that the amounts of data nodes matched to each sub sink grid are almost the same. Besides, the total energy consumption must be minimized simultaneously. The energy-efficiency issue in CAMD is named after Optimal Grid Assignment (OGA) in the remaining content.

Matrix  $M_{(R_g-1)C_g \times C_g}$  is introduced to indicate the relationship between two kinds of grids. Matrix  $H_{(R_g-1)C_g \times C_g}$  records the accumulated shortest hops in each data node grid. In addition, row vector  $V_{1 \times (R_g-1)C_g}$  is adopted to record the amounts of data nodes in each data node grid.

Matrix  $M_{(R_g-1)C_g \times C_g}$  is constructed by elements  $m_{ij}$ , where  $i \in [1, (R_g - 1) \cdot C_g]$  is the sequence number of data node grids and  $j \in [1, C_g]$  represents the sequence number of sub sink grids. The value of  $m_{ij}$  is binary with 1 and 0, and subject to:

$$\sum_{j=1}^{C_g} m_{ij} = 1, \quad \forall i \in [1, (R_g - 1) \cdot C_g] \quad (9)$$

because each data node grid selects but only one sub sink grid as its destination.

Specially, the elements  $h_{ij}$  in matrix  $H_{(R_g-1)C_g \times C_g}$  is the accumulated value of shortest hops from data nodes in the  $i$  th data node grid to the  $j$  th sub sink grid, so it can be gotten by:

$$h_{ij} = \sum_{k=1}^{n_i} h_k^{ij} \quad (10)$$

where  $h_k^{ij}$  represents the shortest hops from node  $i$  to the sub sink grid  $j$ .

Vector  $V_{1 \times (R_g - 1)C_g}$  has  $(R_g - 1) \cdot C_g$  elements, the values of which indicate the amounts of data nodes in each data node grid orderly.

Now, the OGA issue can be described by matrixes discussed above. Firstly, another row matrix  $W_{1 \times C_g}$  that indicates the sum of data nodes matched to each sub sink grid is gotten by  $V_{1 \times (R_g - 1)C_g} \cdot M_{(R_g - 1)C_g \times C_g}$ . Then the constraint function Equation (8) is equivalent to:

$$\min \frac{1}{C_g} \sum_{i=1}^{C_g} (w_{1i} - \bar{w})^2 \quad (11)$$

where  $\bar{w} = \frac{1}{C_g} \sum_{i=1}^{C_g} w_{1i}$  is the arithmetic mean value of all elements  $w_{1i}$  in row matrix  $W_{1 \times C_g}$ .

Then, the objective function Equation (6) is equivalent to:

$$\min \sum_{i=1}^{(R_g - 1)C_g} \sum_{j=1}^{C_g} m_{ij} \cdot h_{ij} \quad (12)$$

where all constants in Equation (6) are ignored because they exert no impacts on the final results.

Finally, the OGA issue is described by objective function Equation (12) as well as constraint functions Equation (9) and Equation (11). The problem is NP-hard (Oncan, 2007) with combinatorial optimization. Genetic algorithm (GA) is an effective method to solve the OGA problem. GA usually has four phases: encoding, selection, crossover and mutation.

### Population Initialization and Selection

A binary chromosome is needed to characterize an individual and the specific chromosome can also be a sample that represents the assignment of all data node grids to sub sink grids. So, the matrix  $M_{(R_g - 1)C_g \times C_g}$  can directly be regarded as the gene code of a chromosome.

The fitness value  $f$  and unfitness value  $uf$  of solutions are formulized by:

$$f = \sum_{i=1}^{(R_g - 1)C_g} \sum_{j=1}^{C_g} m_{ij} \cdot h_{ij} \quad (13)$$

$$uf = \frac{1}{C_g} \sum_{i=1}^{C_g} (w_{1i} - \bar{w})^2 \quad (14)$$

Then  $n_{ga}$  individuals are generated to build the initial population. The generation of the initial population is randomly performed, but it has to obey constraint Equation (9). Besides, the initialization must obey the rule: each data node grid selects a destination sub sink grid randomly among the sub sink grids which are available for it. This rule can avoid infeasible solutions.

Due to the introduction of graphing, the complexity of OGA matching algorithm is decreased. The individuals with the 2nd and 3rd smallest fitness value can directly be selected as parents. This selection algorithm is able to protect the best individual in current population.

### Crossover and Mutation

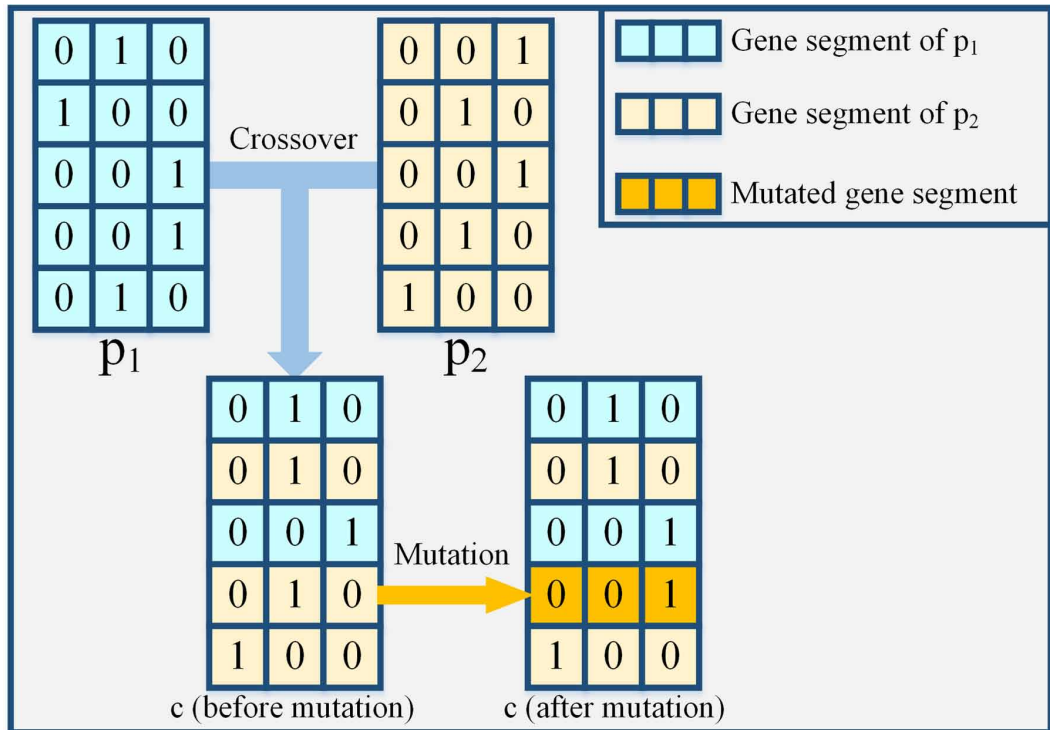
Each row of the chromosome matrix is regarded as a gene segment and all gene segments represent the assignment of all data node grids orderly. To avoid infeasible gene segment of the child solution, parents just exchange their gene segments of the same data node grid.  $p_1$ ,  $p_2$  and  $c$  are used to denote the parents and their child.

The crossover algorithm is based on the fitness value of parents. Firstly, the ratios of each fitness to the sum of the two parents are calculated. Secondly,  $n_{p_1}^{cs}$  gene segments of  $p_1$  are randomly located and substitute corresponding ones of  $p_2$  to produce  $c$ . Besides,  $n_{p_1}^{cs}$  and  $n_{p_2}^{cs}$  are subject to:

$$n_{p_1}^{cs} + n_{p_2}^{cs} = (R_g - 1) \cdot C_g \quad (15)$$

$$\frac{n_{p_1}^{cs}}{n_{p_2}^{cs}} \approx \frac{f(p_2)}{f(p_1)} \quad (16)$$

Figure 5. The process of crossover and mutation



Equations (15) and (16) ensure that the child inherits more genetic information from the parent with smaller fitness. Figure 5 illustrates the crossover algorithm proposed above.

When the crossover operation is finished, there is a little chance for each child to mutate. The mutation mechanism is capable of avoiding premature convergence. If gene mutation occurs in a child solution, a gene segment will be randomly located and the bits of which will be changed according to the same rule on population initialization.

### Population Updating

When the crossover and mutation operation in an iteration is finished, the population may be updated with the child solution.

First of all, the chromosome of child solution is examined. If all bits of the chromosome are the same with any individual in the parent population, the child solution is regarded as identical and will be discarded.

Once a unique child solution is reserved, the fitness value and unfitness value of which are compared with the parent population. If its unfitness value is smaller than the individual with the largest unfitness among the parent population, the child solution will replace this largest individual and the population is updated. Otherwise, if the child solution's unfitness value is equal to the largest unfitness among the parent population, the child solution will replace the largest individual as long as the child's fitness value is smaller than this parent individual. The update algorithm prevents the GA from premature convergence for unfitness. In other conditions, the population keeps unchanged.

### Malicious Node Detection and Forensics

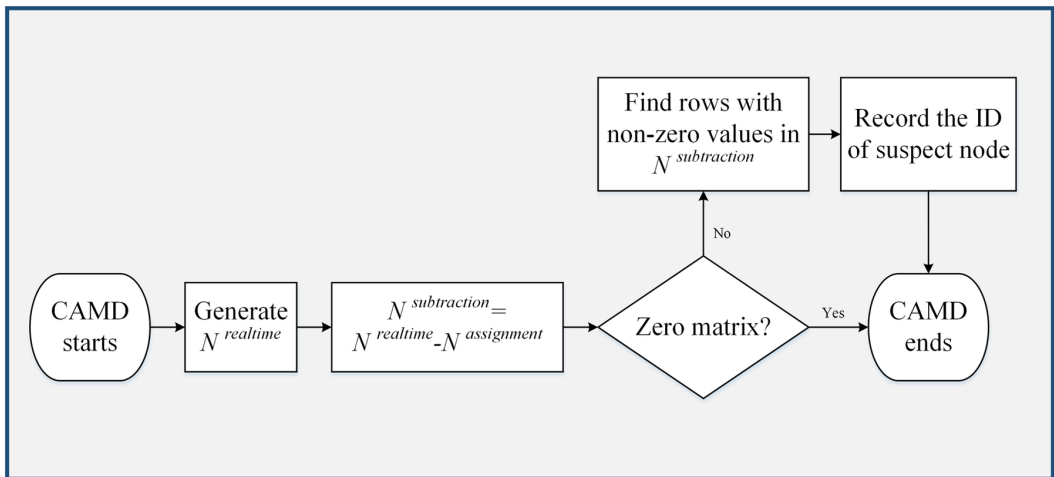
As mentioned above, all data nodes including sub sinks are homogeneous and will generate data from their environment or objects attached. If some data nodes are unfortunately accessed and manipulated by cyber criminals, they will generate fake information and send it to a random sub sink grid. Malicious nodes consume the resources of the whole network and damage the validity of the information gathered.

When the deployment and graphing of IoT are finished, the information of the network system can be described by the matrixes proposed in preceding sections. Based on the grid assignment, a Centralized and Active Malicious Node Detection (CAMD) with information comparison mechanism is designed and installed in the mobile sink to detect the malicious node manipulated by cyber criminals and to provide strong digital evidence for forensics. To locate the malicious nodes, CAMD checks the assignment of the grids. The matrix-based information comparison mechanism in CAMD is illustrated by Figure 6.

When the mobile sink arrives at either end point of its path, which signifies that all buffered data in sub sinks has been gathered by the mobile sink, the mobile sink will generate a binary matrix denoted by  $N_{n \times C_g}^{realtime}$ . The value of a certain element  $n_{ij}^{realtime}$  in  $N_{n \times C_g}^{realtime}$  indicates that whether the data packet of a certain data node  $i$  arrives at sub sink grid  $g_j^{ss}$  finally. Then, the mobile sink compares  $N_{n \times C_g}^{realtime}$  to the assignment matrix  $N_{n \times C_g}^{assignment}$  which is generated by the genetic algorithm. The comparison is conducted by matrix subtraction. If the result matrix of the subtraction, denoted by  $N_{n \times C_g}^{subtraction}$  is not a zero matrix, malicious nodes are believed to exist in the network. What follows next is to check every row of  $N_{n \times C_g}^{subtraction}$  orderly and the one with non-zero value represents a malicious node.

According to the result of the matrix comparison mechanism of CAMD, the location and identity of malicious nodes that are manipulated by cyber criminals can be precisely detected. Then the information gathered by the mobile sink which is from the detected malicious nodes can be regarded

Figure 6. CAMD algorithm illustration



as fake or harmful information. The binary subtraction matrix  $N_{n \times C_g}^{subtraction}$  records the fake information reporting behavior conducted by cyber criminals and the digital result is strong and reliable evidence for forensics.

## DISTRIBUTED AND PASSIVE EEA RESISTANCE (DPER)

If cyber criminals make the manipulated nodes transmit data packet more frequently, the nodes along the data packet delivery path will deplete more rapidly. In this section, two communication protocols are devised for DPER to alleviate the impact of such EEA attack. The protocols are installed in each data node, and a grid-based routing table is installed in each node as shown in Table 2.

In the grid routing table, the pipe nodes as well as the grids they belong to are recorded. In the grid routing table, there is a pointer named after *current\_pipe*. The data packet forwarded by this data node will be transmitted to the pipe node which is pointed to by *current\_pipe*. Once the packet is transmitted, *current\_pipe* will point to the next pipe node circularly. With the grid routing table, the

Table 2. Grid routing table

ID of Local Node	ID of Sub Sink Grid	Hops of Shortest Path	ID of Gateway Grid	ID of Accessible Pipe Nodes
28	1	5	11	* 7
				15
				45
	2	6	11	* 7
				15
				45
	5	7	15	* 5
				18

communication protocols for DPER are capable of dynamically balancing the energy consumption of data nodes during the data gathering period and preventing some node from depletion caused by EEA.

The communication protocols are made up of network discovery and packet delivery. The protocols are designed based on the graph and Shortest Path Tree (SPT), therefore the technique can be regarded as Advanced Shortest Path Tree (ASPT).

### Graph-based Network Discovery Protocol

The main task of this protocol is to lay a basis for the data packet deliver process by establishing the connection among grids. To achieve these goals, the mobile sink needs to accomplish three round trips. The grid routing table and the assignment result of the genetic algorithm are also broadcasted to each node by this protocol.

For the first trip, the mobile sink broadcasts sub sink checking packets at a time interval  $t_{int}^{ck}$ . Once a certain data node receives a sub sink checking packet, it becomes a sub sink logically and begin the SPT establishment on its own. When all SPT setup packets are transmitted, each data node record the shortest paths and gateways from itself to all reachable sub sinks. Then data node sends an ordinary data packet to the nearest sub sink with the information of its location and SPT routing information. If the data packets are successfully transmitted, data nodes will clear its memory for grid routing information which is established in future.

During the second round trip, the mobile sink will broadcast information request packets, and the sub sinks will reply to the mobile sink with a data packet containing the information of data nodes collected during the first trip. If the data packets are transmitted successfully, sub sinks will delete the information. When the mobile sink finishes its second round trip, it will perform the graphing operation and GA-based matching algorithm.

The task during the third round trip is to inform data nodes the result of the GA-based matching algorithm and establish the ASPT. The format of the packet during this period is defined as  $OGA\_Pkt \{g\_hop, src, gtw\_id, pipe\_adr, mapping\_list\}$ , where  $g\_hop$  is the accumulated hop from the sub sink grid  $src$  to the grid that contains the node receive the packet.  $gtw\_id$  and  $pipe\_adr$  respectively represent identification of the grid this node belongs to and the address of pipe node which forwards this packet.  $mapping\_list$  contains information including the algorithm results. The mobile sink broadcasts  $OGA\_Pkt \{0, 0, 0, 0, mapping\_list\}$  during the trip. Data nodes which receive  $OGA\_Pkt$  acquire the identification of its destination sub sink grid  $g_i^{ss}$  as well as the identification of grid to which it belongs. Meanwhile, data node will update its grid routing table with smallest hops from itself to reachable sub sink grids. Before forwarding the OGA packet, the node deletes its own item from  $mapping\_list$ .

### Dynamic Adaptive Packet Delivery Protocol

EEA will boost the data nodes' energy consuming speed. During the detection delay caused by centralized IDS mechanism, the influence of EEA must be alleviated. Conventional SPT protocol provide only one gateway node to each data node. Once the routing table for this node is established, all packets in the future from this node will be forwarded by the only gateway node provided by SPT protocol. The gateway node will deplete due to heavy packet traffic. If EEA attacks are conducted by cyber criminals, the fast energy dissipation situation will be deteriorated.

Benefitted by the operation of graphing, all pipe nodes which are potential gateway nodes found by the network discover protocol. With these pipe nodes, dynamic gateway method can be easily implemented.

During the data delivery phase, each data node consults its grid routing table to find out the pipe node by which its data packet is going to be forwarded. Then the data node transmits the data packet to  $current\_pipe$ . The format of data packets is defined as  $Data\_Pkt \{src, dst, pipe, buffer\}$ .

When a certain node receives a *Data\_Pkt*, it firstly compares the destination sub sink grid of this packet with the node's own grid. If the receiving node happens to belong to the destination sub sink grid of this packet, which means the data packet reaches its destination, the receiving node (sub sink) stores the information of this data packet in its memory and waits to send it to the mobile sink. Otherwise, if the receiving node is not a sub sink, it will forward the data packet to *current\_pipe*. Then the *current\_pipe* points to the next pipe node.

## SIMULATION IMPLEMENTATION AND RESULTS ANALYSIS

The performance of the hybrid IDS proposed was evaluated on the Network Simulator-3 (NS-3, version 3.25) platform (Lacage, Carneiro, & et al., 2008). In the simulation experiments, the area of  $A$  is  $400m \times 600m$ . The speed of the mobile sink was  $v_m = 5m / s$ . The data generation rate was  $p = 200bits / s$  and the data transmission interval  $t_{int}$ . The initial energy of each data node was  $E_{int} = 10Joules$  and the energy consumption factor was  $e = 0.5\mu J / bit$ . The maximum communication radius was  $r = 50\sqrt{2}m$  and the length of grid side was  $50m$ . The data transmission rate between any two nodes was  $44Kb / s$ . The possibility that data nodes might be manipulated by cyber criminals during each data interval was 2%.

There were 14 network scales which increased from 110 nodes to 240 nodes. Combined with the size of the region  $A$  and communication radius  $r$ , the network scales varied from low-density distribution to high-density distribution.

According to the previous discussion, the following metrics were adopted to evaluate the hybrid IDS and data gathering scheme:

- *Network Lifetime* was defined as the time duration from the beginning of mobile sink assignment to the first energy exhaustion of any sink node or the dynamic data gathering percentage drops below the threshold. The network lifetime represented the energy-efficiency of the hybrid IDS including DPER.
- *Network Energy Consumption Efficiency* (NECE) was defined to measure the power of the whole network and NECE was the ratio of total energy consumed by all data nodes to the network lifetime.
- *Malicious Node Detection Ratio* (MNDR) represented the performance of matrix comparison mechanism in CAMD.
- *Futile Data Generation* represented the impact of the hybrid IDS on decreasing the futile data generated by malicious nodes.

### Impact of the Hybrid IDS on Network Lifetime

Due to the hierarchy of the topology, it was quite reasonable to treat the lifetime of the first exhausted sub sink as the lifetime of the whole network. However, the definition might cause a new problem: a lot of ordinary data nodes depleted before sub sinks and the data gathering performance would decline. Thus, a threshold on dynamic data gathering percentage was introduced to avoid the extreme situation. The threshold was set to be 90% to ensure that the data gathering performance was acceptable.

It can be seen from Figure 7 that the network lifetime was terribly influenced by the EEA conducted by the malicious nodes. EEA would deplete data nodes' energy sources more quickly than normal nodes. What's worse is that when the manipulated nodes exhausted, the data delivery performance deteriorated and eventually the network lifetime came to an end ahead of schedule. With the hybrid IDS proposed, the malicious nodes manipulated by cyber criminals could be detected precisely by CAMD and the influence of the EEA were reduced accordingly. Particularly, DPER

balanced the energy consumption among pipe nodes. Besides the manipulated nodes conserved energy passively but effectively.

Moreover, the GA-based data gathering scheme was compared with traditional SPT and RANDOM schemes. The results were illustrated in Figure 8. The simulation results indicated that the GA-based data gathering scheme proposed showed superiority.

### Impact of the Hybrid IDS on Network Energy Consumption Efficiency (NECE)

The Network Energy Consumption Efficiency (NECE) represented the accumulated power of all the nodes in the IoT network. NECE was affected the network lifetime to some extent. Before the malicious nodes were detected by CAMD, cyber criminals would perform EEA attacks and the energy consumed by victim nodes greatly increased. If the detection took too much time due to the centralized method, the influence of EEA would go on and the NECE level would increase.

Figure 9 showed that the NECE level got much higher when malicious nodes existed and EEA was conducted. However, with the hybrid IDS proposed, malicious nodes would be detected and corrected in time. The impact of EEA got reduced accordingly. Distributed methods might outperform CAMD in the form of real-time performance, but the extra energy consumption of the whole network was unacceptable to resource-constrained IoT applications. DPER diluted the impact of EEA and the impact of the detection delay was alleviated. The associated work of DPER and CAMD improved the Network Lifetime of IoT.

Figure 7. The impact of the hybrid IDS proposed on network lifetime

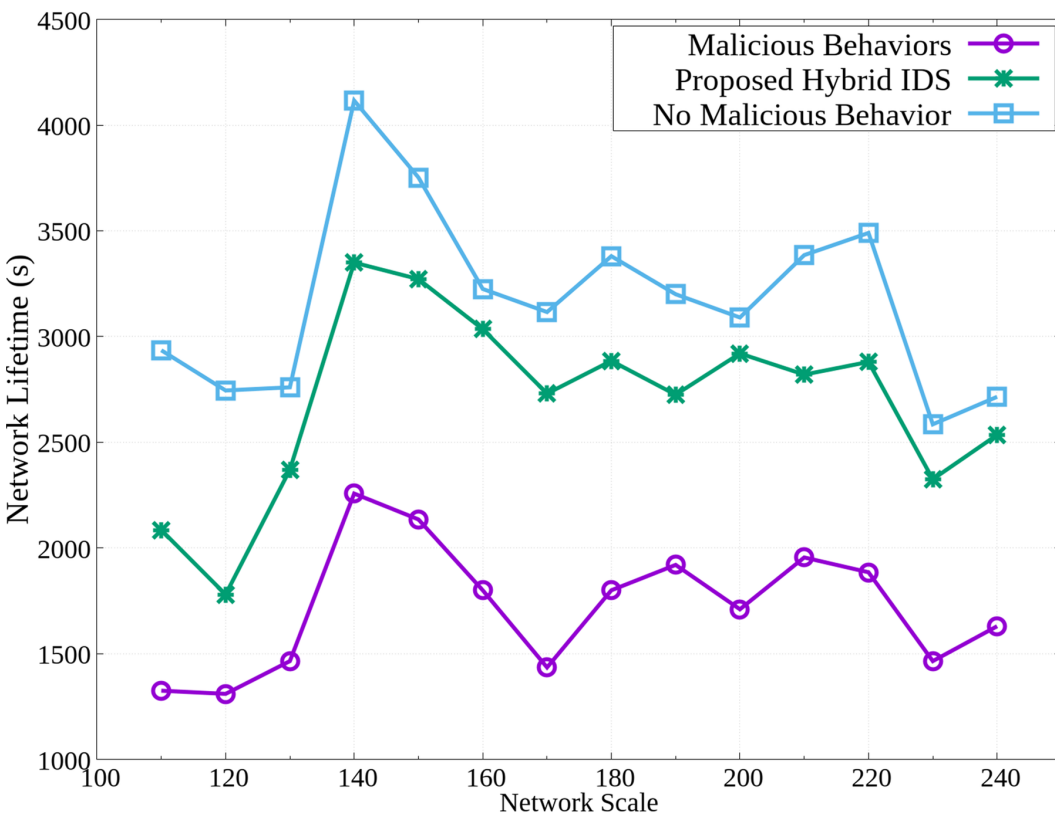
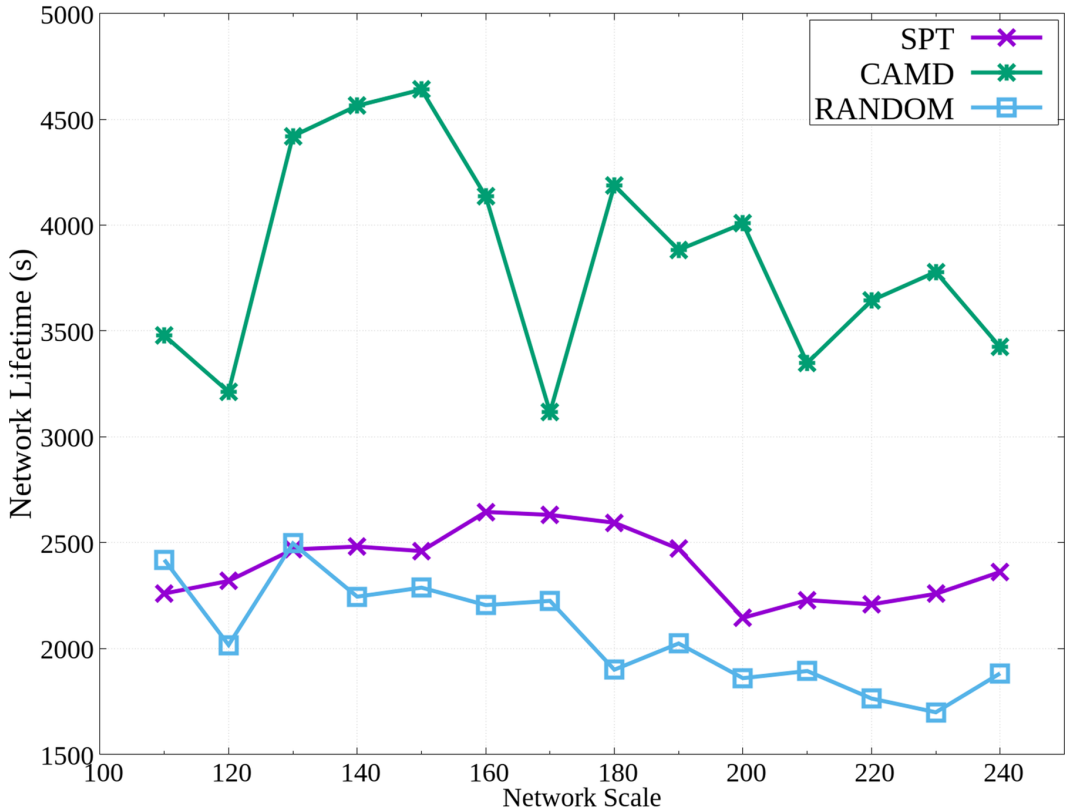




Figure 8. Network Lifetime results compared with SPT and RANDOM schemes



### Impact of CAMD on Malicious Node Detection Ratio (MNDR)

There was a little possibility that the malicious node completely copied the behaviors of the normal node, which meant the destination sub sink grid randomly selected by the malicious node happened to be the one OGA algorithm assigned. The proposed CAMD could not deal with this situation so far, because CAMD was a centralized method. However, the probability was very low in large-scaled deployments as shown in Figure 10.

### Impact on Futile Data Generation

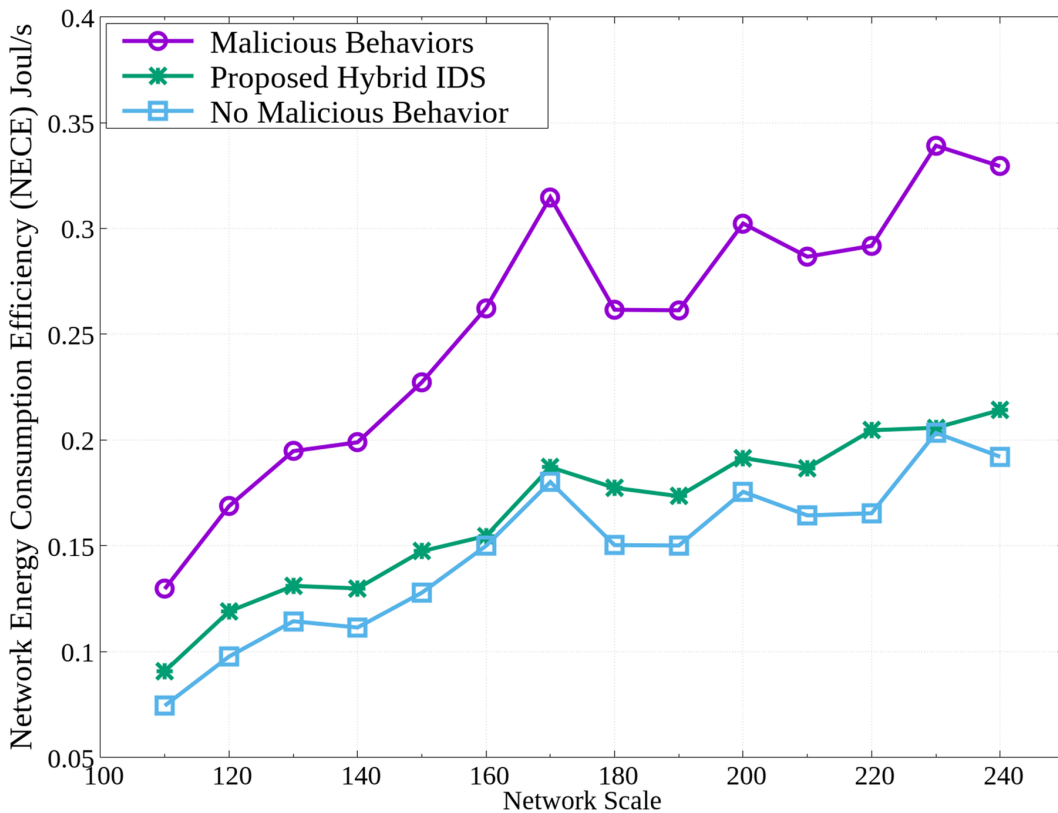
The mobile sink identified and located malicious nodes by CAMD, then the maintainer of the network would get notification and the malicious nodes in the network would get fixed or replaced in a short time. Then the nodes would no longer act malicious behaviors.

The fake data generated by cyber criminals could be regarded as futile data. Figure 11 demonstrated the total amount of futile data caused by malicious nodes. Without the hybrid IDS, the amount of futile data gathered was quite high and the hybrid IDS effectively reduced the harmful data.

## CONCLUSION

IoT applications with path-constrained mobile sink can be applied in various practical scenarios to achieve information collection and environment monitoring. IoT networks are always privacy sensitive and even confidential. However, resource constraint makes traditional network security and network forensics technologies un-appropriate for IoT.

Figure 9. The impact on NECE of the hybrid IDS



Based on graph theory a hybrid IDS is proposed. A Centralized and Active Malicious Detection (CAMD) method is integrated in the Genetic Algorithm-based (GA) data gathering scheme. CAMD detects malicious nodes manipulated by cyber criminals and provides strong digital evidence for forensics. Then Distributed and Passive EEA Resistance (DPER) is implemented through a set of advanced Shortest Path Tree-based (ASPT) communication protocols to alleviate the impact of Energy Exhaustion Attacks (EEAs) conducted by cyber criminals.

Simulation results performed on NS-3 showed that the hybrid IDS performed as anticipated. The cyber crimes in the form of fake information reporting and EEA attacks were detected and alleviated. Reliable digital evidence was provided by the hybrid IDS. Besides, the energy efficiency of the IoT network was not deteriorated by the hybrid IDS.

## ACKNOWLEDGMENT

This research was supported by the YangFan Innovative & Entrepreneurial Research Team Project of Guangdong Province; the National Natural Science Foundation of China [grant number 61502050]; the Beijing Key Laboratory of Work Safety Intelligent Monitoring; the Key Research and Develop Project of Artificial Intelligence Technology Innovation Major Thematic Projects of Chongqing [grant number cstc2017rgzn-zdyfx0003]; the National Key Research and Development Project [grant number 2018YFB1600800].

Figure 10. Impact on malicious node detection ratio (MNDR)

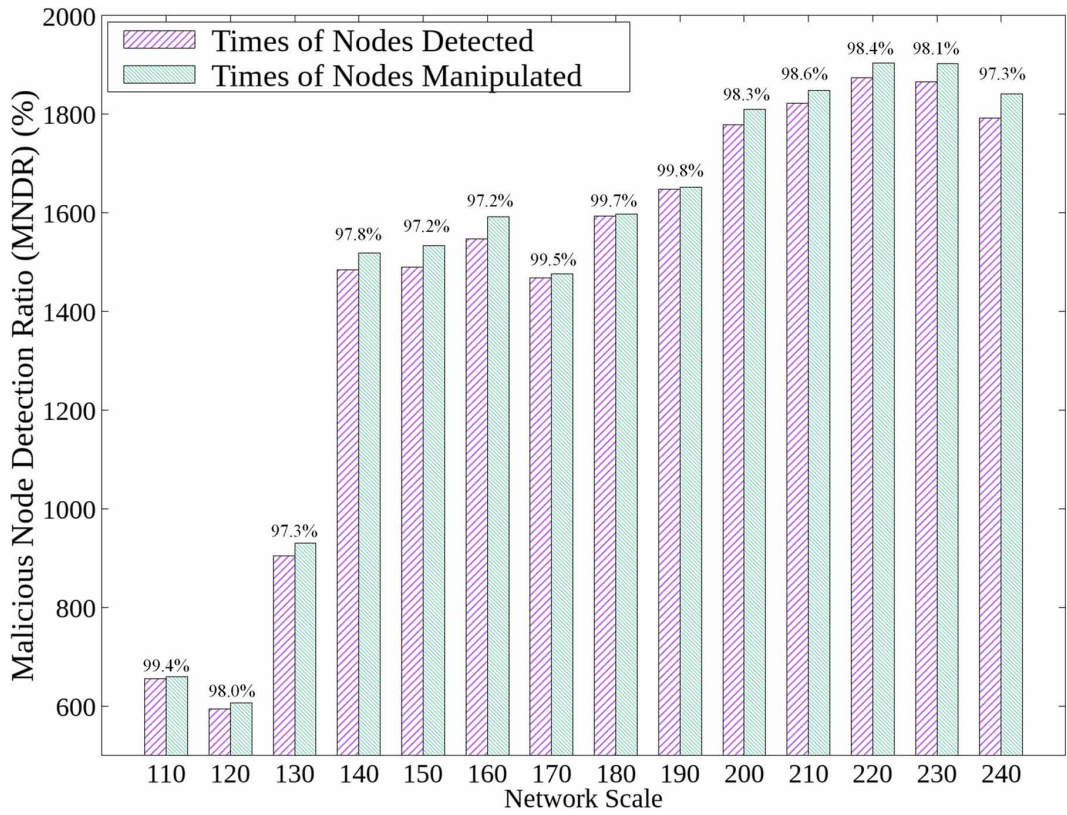
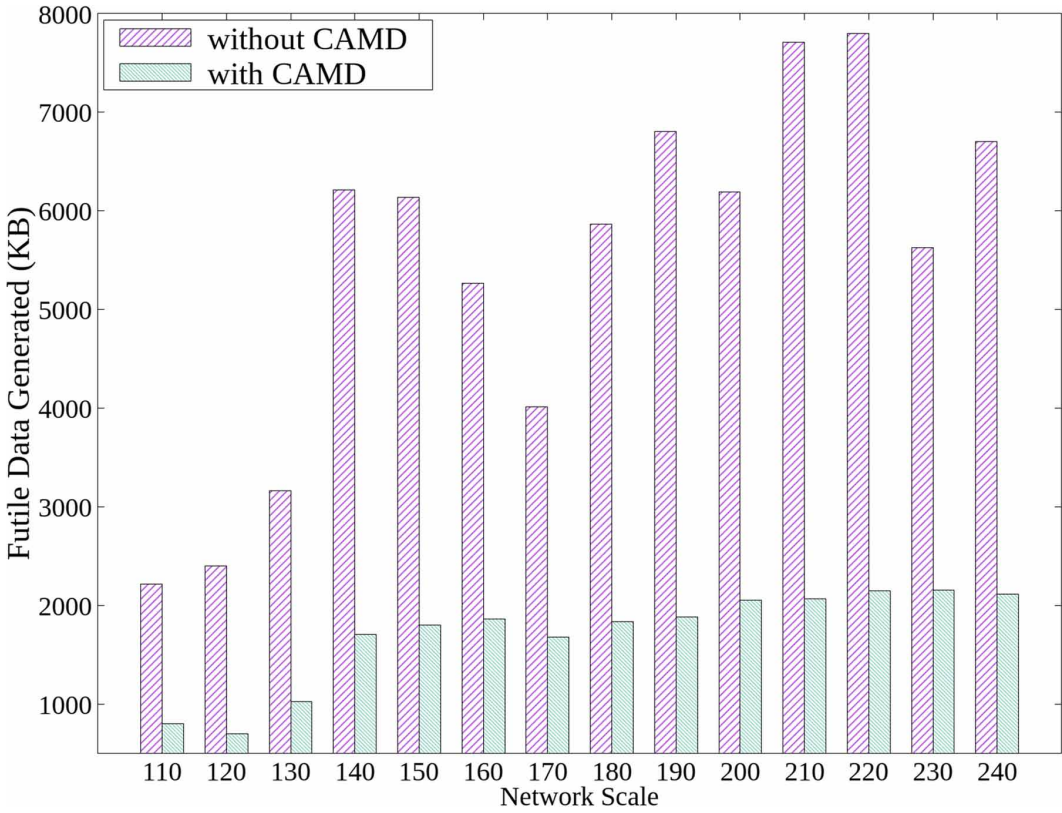


Figure 11. The impact on futile data generation



## REFERENCES

- Adeyemi, I. R., Razak, S. A., & Nor Azhan, N. (2013). A Review of Current Research in Network Forensic Analysis. *International Journal of Digital Crime and Forensics*, 5(1), 1–26. doi:10.4018/jdcf.2013010101
- Ahmed, A. A. (2017). Investigation Approach for Network Attack Intention Recognition. *International Journal of Digital Crime and Forensics*, 9(1), 17–38. doi:10.4018/IJDCF.2017010102
- Alaba, F. A., Othman, M., Hashem, I. A. T., & Alotaibi, F. (2017). Internet of Things security: A survey. *Journal of Network and Computer Applications*, 88, 10–28. doi:10.1016/j.jnca.2017.04.002
- Alrajeh, N. A., Khan, S., Lloret, J., & Loo, J. (2014). Artificial Neural Network Based Detection of Energy Exhaustion Attacks in Wireless Sensor Networks Capable of Energy Harvesting. *Ad-Hoc & Sensor Wireless Networks*, 22(1-2), 109–133.
- Amaral, J., Oliveira, L., Rodrigues, J., Han, G., & Shu, L. (2014). Policy and Network-based Intrusion Detection System for IPv6-enabled Wireless Sensor Networks. In *Proceedings of International Conference on Communications, ICC 2014* (pp. 1796-1801). IEEE. doi:10.1109/ICC.2014.6883583
- Gunasekaran, M., & Periakaruppan, S. (2017). GA-DoSLD: Genetic Algorithm Based Denial-of-Sleep Attack Detection in WSN. *Security and Communication Networks*, 2017, 1–10. doi:10.1155/2017/9863032
- Hasan, M. Z., Al-Rizzo, H., & Günay, M. (2017). AI-Rizzo, H., & Gunay, M. (2017). Lifetime Maximization by Partitioning Approach in Wireless Sensor Networks. *EURASIP Journal on Wireless Communications and Networking*, 2017(1), 15. doi:10.1186/s13638-016-0803-1
- Huang, H. L., & Savkin, A. V. (2016). Optimal Path Planning for a Vehicle Collecting Data in a Wireless Sensor Network. In *Proceedings of 35th Chinese Control Conference (CCC)* (pp. 8460-8463). Academic Press. doi:10.1109/ChiCC.2016.7554706
- Joby, P. P., & Sengottuvelan, P. (2015). A Localised Clustering Scheme to Detect Attacks in Wireless Sensor Network. *International Journal of Electronic Security and Digital Forensics*, 7(3), 211–222. doi:10.1504/IJESDF.2015.070386
- Lacage, M., & Carneiro, G. (2008). *Network Simulator-3*. Retrieved from <https://www.nsnam.org>
- Le, A., Loo, J., Chai, K.K., & Aiash, M. (2016). A Specification-based IDS for Detecting Attacks on RPL-based Network Topology. *Information*, 7(2).
- Lee, T. H., Wen, C. H., Chang, L. H., Chiang, H. S., & Hsieh, M. C. (2014). A Lightweight Intrusion Detection Scheme Based on Energy Consumption Analysis in 6LowPAN. In *Proceedings of Advanced Technologies, Embedded and Multimedia for Human-Centric Computing, HumanCom and EMC 2013* (pp. 1205-1213). Academic Press.
- Liu, Q., Zhang, K., Shen, J., Fu, Z., & Linge, N. (2016). Glrm: An Improved Grid-based Load-Balanced Routing Method for WSN with Single Controlled Mobile Sink. In *Proceedings of 18th International Conference on Advanced Communications Technology: "Information and Communications for Safe and Secure Life!" ICACT 2016* (pp. 34-38). Academic Press.
- Nake, N. B., & Chatur, P. N. (2016). An Energy Efficient Grid Based Routing in Mobile Sink Based Wireless Sensor Networks. In *Proceedings of 2nd International Conference on Science Technology Engineering and Management, ICONSTEM 2016*. Academic Press. doi:10.1109/ICONSTEM.2016.7560872
- Oh, D., Kim, D., & Ro, W. W. (2014). A Malicious Pattern Detection Engine for Embedded Security Systems in The Internet of Things. *Sensors (Basel)*, 14(12), 24188–24211. doi:10.3390/s141224188 PMID:25521382
- Oncan, T. (2007). A Survey of the Generalized Assignment Problem and its Applications. *Information Systems and Operational Research*, 45(3), 123–142. doi:10.3138/infor.45.3.123
- Qin, W., Hempstead, M., & Yang, W. (2006). A Realistic Power Consumption Model for Wireless Sensor Network Devices. In *Proceedings of 3rd Annual IEEE Communications Society on Sensor and Ad Hoc Communications and Networks* (pp. 286–295). SECON.
- Schwartz, E. (2010). Network packet forensics. In Mozayani & Ashraf (Eds.), *Cyber Forensics: Springer's Forensic Laboratory Science Series* (pp. 85-101). Berlin, Germany: Springer. doi:10.1007/978-1-60761-772-3\_7

Smeets, H., Shih, C. Y., Zuniga, M., Hagemeyer, T., & Marrón, P. J. (2013). Trainsense: A Novel Infrastructure to Support Mobility in Wireless Sensor Networks. In *Proceedings of Wireless Sensor Networks - 10<sup>th</sup> European Conference, EWSN 2013* (pp. 18-33). Academic Press. doi:10.1007/978-3-642-36672-7\_2

Teklemariam, G. K., Van Den Abeele, F., Moerman, I., Demeester, P., & Hoebeke, J. (2016). Bindings and RESTlets: A Novel Set of CoAP-Based Application Enablers to Build IoT Applications. *Sensors (Basel)*, 16(8), 1217. doi:10.3390/s16081217 PMID:27490554

Wallgren, L., Raza, S., & Voigt, T. (2013). Routing Attacks and Countermeasures in The RPL-based Internet of Things. *International Journal of Distributed Sensor Networks*, 2, 167–174.

Yun, Y. S., Xia, Y., Behdani, B., & Smith, J. C. (2013). Distributed Algorithm for Lifetime Maximization in a Delay-tolerant Wireless Sensor Network with a Mobile Sink. *IEEE Transactions on Mobile Computing*, 12(10), 1920–1930. doi:10.1109/TMC.2012.152

Zarpelao, B. B., Miani, R. S., Kawakani, C. T., & de Alvarenga, S. C. (2017). A survey of intrusion detection in Internet of Things. *Journal of Network and Computer Applications*, 84, 25–37. doi:10.1016/j.jnca.2017.02.009

Zhang, Y. Q., Zhou, Z. B., Zhao, D., Barhamgi, M., & Rahman, T. (2017). Graph-based mechanism for scheduling Mobile Sensors in Time-sensitive WSNs Applications. *IEEE Access*, 5, 1559–1569. doi:10.1109/ACCESS.2017.2667687

Zhou, Z., Du, C., Shu, L., Hancke, G., Niu, J., & Ning, H. (2016). An Energy Balanced Heuristic for Mobile Sink Scheduling in Hybrid WSNs. *IEEE Transactions on Industrial Informatics*, 12(1), 28–40. doi:10.1109/TII.2015.2489160

*Chao Wu received the B.S. degree in Elec-tronic Information Science and Technology from Chongqing University, China in 2009. He received the Ph.D. degree with the School of Electronic Engineering, Beijing University of Posts and Telecommunications, Beijing. He engaged in the research of Wireless Sensor Networks and Internet of Things. Presently, he is with Chongqing Vehicle Test & Research Institute Co. Ltd., Chongqing.*

*Yuan'an Liu received the M.E. and Ph.D. degrees in Chengdu University of Electronic Science and Technology, China, in 1989 and 1992. He engaged in the post-doctoral research in Beijing University of Posts and Telecommunications (BUPT), China, from 1992 to 1994. He worked at the Carleton University, Canada, during 1995-1997. Now he is the executive director of School of Electronic Engineering, BUPT, majoring in EMC, mobile communications, and Internet of Things. He is a fellow of IEE. He joined the 26th Institute of Electronic Ministry of China, Langfang, China, to develop the inertia navigating system in 1984. In 1992, he held a first post-doctoral position with the Electromagnetic Compatibility (EMC) Laboratory, BUPT, Beijing, China, and a second post-doctoral position with the Broadband Mobile Laboratory, Department of System and Computer Engineering, Carleton University, Ottawa, ON, Canada, in 1995. Since 1997, he has been with the Wireless Communication Center, College of Telecommunication Engineering, BUPT, as a Professor, where he is involved in the development of next-generation cellular system, wireless LAN, Bluetooth applications for data transmission, EMC design strategies for high speed digital system, and electromagnetic interference and electromagnetic susceptibility measuring sites with low cost and high performance.*

*Fan Wu received her Ph.D. degree from Beijing University of Posts and Telecommunications, China, in 2009. Now she is an associate professor at School of Electronic Engineering, Beijing University of Post and Telecommunication. Her research interests include Internet of Things, sensor search, and pervasive computing.*

*Feng Liu is with State Key Laboratory of Information Security, Institute of Information Engineering and School of Cybersecurity University of Chinese Academy of Sciences, Chinese Academy of Sciences.*

*Wenhao Fan received his Ph.D. degree in Beijing University of Posts and Telecommunications, China, in 2013. Now he is a master tutor at School of Electronic Engineering, Beijing University of Post and Telecommunication. His research interests include mobile devices and cloud computing.*

*Bihua Tang received the M.S. degree in Chengdu University of Electronic Science and Technology, China, in 1989. She worked at the Carleton University, Canada, in 1996. She spent one year as a senior visiting scholar in National Institute of Informatics, Japan, in 2002. Now, she is a professor at School of Electronic Engineering, Beijing University of Post and Telecommunication, China. Her research interests include Internet of Things, mobile computing, and wireless sensor networks.*