

# On Data Aggregation with Secure Bloom Filter in Wireless Sensor Networks

Zhijun Li

Department of Electrical and Computer Engineering  
University of Waterloo, Waterloo, Ontario, Canada  
Email: leezj@engmail.uwaterloo.ca

Guang Gong

Department of Electrical and Computer Engineering  
University of Waterloo, Waterloo, Ontario, Canada  
Email: ggong@calliope.uwaterloo.ca

## Abstract

Designing secure data aggregation schemes, which are critical to many sensor network applications, imposes interesting and formidable challenges. In this paper, we propose a succinct and practical secure aggregation protocol by combining HMAC (associated with a cryptographic hash function) with Bloom filter, which then is defined as secure Bloom filter. Unlike most previous approaches, which are aimed to provide security mechanisms for ordinary aggregation operations, our proposal firstly is an effective aggregation protocol, suitable for a specific but popular class of aggregation in wireless sensor networks. Benefiting from secure Bloom filter, the protocol, without any unrealistic assumptions, fulfills the fundamental security objective of preventing outside adversaries and compromised inside nodes from harming the overall network result. We systematically analyze the protocol performances and run extensive simulations on different network scenarios for evaluation. The simulation results demonstrate that the proposed protocol presents remarkable performances on security, communication cost, and degree of node energy consumption balance.

## I. INTRODUCTION

In typical wireless sensor networks (WSNs), hundreds and thousands of low-cost sensor nodes scatter within a surveillance area, receive commands from a base station, perform designated detection tasks accordingly, and collaboratively transmit results back to the base station in a multi-hop, infrastructure-less communication pattern. In many cases, instead of forwarding every individual message to the base station, sensor network protocols support data aggregation—the operation by intermediate nodes that combines many messages and sends out aggregated results. As a matter of fact, from the very beginning of wireless sensor networks development [1], [2], [3], [4], [5], it has already been widely accepted that data aggregation plays a critical role in the practicability and appealing of WSNs. Due to infeasibility of recharging nodes batteries in most circumstances, energy becomes the most valuable resource for sensor nodes. Among all nodes operations, data transmission consumes energy the most [2], [5]. Intuitively, data aggregation during message transmission is an effective method to preserve sensor nodes precious energy. Moreover, in the absence of data aggregation, sensor nodes near the base station will suffer from heavy message transmission overhead, and then die of power exhaustion much sooner than other nodes, breaking down the whole network's functionality. As a result, data aggregation attracts a great deal of attention and many a data aggregation scheme has been proposed in recent years. Interested readers may refer to [3], [4] for systematic surveys on this topic.

When sensor nodes are deployed in a hostile environment, security measurements should be taken into consideration for network protocols. In general, attacks to wireless sensor networks not only come from outside adversaries, but also can be conducted by compromised, previously legitimate nodes. This kind of attack model, along with the fact that sensor nodes are very resource-constrained, imposes formidable challenges for designing secure data aggregation. When purely cryptographic mechanisms are used for securing data aggregation, homomorphic primitives might be the only suitable candidates, and then many literatures follow this track. However, as we will argue later, the application of homomorphic primitives on data aggregation is highly limited. For other kinds of secure data aggregation approaches, sophisticated protocols are designed to detect malicious behaviors/nodes. Some of those mechanisms rely on unrealistic assumptions, while other involve heavy communication overhead, which conflicts with the very intention of data aggregation and makes it hard to evaluate their applicability. There are some well-designed secure data aggregation schemes, which understandably have different trade-offs and are suitable for particular operations.

In this paper, we propose a secure data aggregation protocol which fits a specific but popular class of aggregations in wireless sensor networks. Unlike most previous approaches, which are aimed to provide security mechanisms for ordinary aggregation operations, our proposal in the first place is an effective and succinct aggregation protocol, which is equipped with built-in security mechanisms, fulfilling the fundamental security purpose of preventing compromised nodes as well as outside adversaries from harming network aggregated results. The security technique backing up our protocol is to combine HMAC, a conventional and provably secure cryptographic primitive, with the popular, space-efficient data structure Bloom filter [6]. The resulting data structure is then defined as secure Bloom filter, and its security implications are addressed in this paper. Then we apply it to construct a secure aggregation protocol for sensor networks. The theoretical analysis and extensive simulation

results demonstrate that the proposed protocol presents remarkable performances on security, communication cost, and energy consumption balance degree among sensor nodes.

The rest of this paper is organized as follows. First, previous work on secure data aggregation is discussed in Section II. Then we state application scenarios, and define network and adversary models in Section III. Afterward, we present secure Bloom filter in Section IV. Section V details our proposed secure aggregation protocol and corresponding analysis. The simulation design and results are elaborated in Section VI. Finally, we conclude our work in Section VII.

## II. PREVIOUS WORK

Homomorphic primitives, besides standard cryptographic functionalities, allow users without secret key to legitimately perform acceptable algebraic operations on protected data block. Since aggregation is essentially some operation, it is intuitive to use homomorphic primitives for securing data aggregation, and a number of approaches [7], [8], [9] use homomorphic encryption for this application. Generally, it is very difficult to design secure symmetric homomorphic block encryption, whereas stream ciphers naturally support homomorphic exclusive-OR operations, which is exactly utilized in [8]. In contrast, public-key homomorphic encryption is an interesting topic, and there exist several relatively practical public-key homomorphic cryptosystems, such as Unpadded RSA, El-Gamal, Goldwasser-Micali, Benaloh and Paillier [10], though all of them only support limited operations on ciphertexts. In 2009, Gentry [11], for the first time, presented a fully homomorphic encryption scheme, which outstandingly allows arbitrary operations on ciphertext. Even though the only two fully homomorphic encryption schemes [11], [12] by now cannot provide competitive performance for most applications, practical fully homomorphic encryption is expected to appear eventually.

Unfortunately, sole public-key homomorphic encryption does not suffice for secure data aggregation in sensor networks, because then anyone can maliciously insert or manipulate results. Even for the stream homomorphic encryption in [8], its security demands message integrity mechanisms as well, and their solution to that issue is essentially a homomorphic MAC. The topic of using homomorphic MAC and homomorphic hashing in secure data aggregation is systematically discussed in [13]. Applying those kinds of homomorphic primitives to data aggregation either involves cumbersome key management problems, or incurs considerable communication costs. Consequently, their applicabilities are highly limited.

Hu and Evans [14] described a secure hop-by-hop data aggregation scheme, and later Jadia and Muthuria [15] extended their scheme. However, both schemes are only capable of preventing a single inside malicious node at an appreciable communication cost, rendering them impractical. Yang et al. [16] presented a secure hop-by-hop data aggregation protocol for sensor networks named SDAP, using the principles of divide-and-conquer and commit-and-attest to design complicated algorithms which cause significant transmission overhead, and may cancel off all communication benefits from data aggregation.

Przydatek, Song, and Perrig [17] proposed secure information aggregation (SIA) to identify forged aggregation values from malicious nodes. In the SIA scheme, a special node named aggregator computes an aggregation result over raw data together with a commitment to the data based on a Merkle-hash tree and sends them back to a remote user, which later challenges the aggregator to verify the aggregation. Later Chan, Perrig, and Song [18] built a hierarchical data aggregation on the aggregate-commit-prove framework in [17], but extended their single aggregator model to a fully distributed setting. Frikken and Dougherty [19] further improved the Chan-Perrig-Song scheme. Moreover, Chan and Perrig [20] derived several security primitives from this kind of algorithms.

## III. BACKGROUND

### A. Application Scenarios

As a typical data aggregation application scenario of our proposal, a wireless sensor network is employed in a hostile environment to detect a number of *predefined, basic* events with only two states: ON and OFF. For most of time, these events stay OFF: they did not occur. When sensor nodes detect an event taking place, they notify the base station about their discovery. Since there are many nodes simultaneously detecting same events, data aggregation is desired to reduce communication costs and to balance nodes energy consumption. The base station may be tolerant of a little inaccuracy on final aggregated results.

Even though the scenario looks like relatively simplistic, it is fairly popular for wireless sensor network applications. For instance, an intrusion-detection sensor network [21] successfully triggers an alarm for a perimeter intrusion event. Moreover, sophisticated tasks can be decomposed into many basic events. In such cases, sensor nodes are responsible for simple event detection, which significantly reduces their production expenditure, while the base station performs advanced data analysis. Other similar application scenarios are voting and counting mechanisms for sensor networks. For example, in many secure sensor network protocols, malicious nodes are determined by majority results from nodes voting; and under many circumstances, the base station may need to frequently count the number of nodes that match some standards.

### B. Network and Adversary Models

We consider an ordinary, homogeneous sensor network that consists of  $n$  low-cost sensor nodes which are resource-constrained in terms of computation, communication, storage, and power supply. Sensor nodes are especially sensitive to energy

consumption and thus data aggregation is desired. As a rule, there exists a central, powerful base station that is responsible for collecting sensor network results. Sensor nodes are not equipped with tamper-proof hardware, and the adversary is capable of compromising and fully controlling arbitrary number of sensor nodes, but we assume that the base station is immune to all physical attacks and it is trustful. Moreover, the adversary can eavesdrop and alter any messages from integrity nodes. In a word, the adversary is granted the full-scale attack capacity against the sensor network except the base station.

In our network model, sensor nodes may be densely deployed, or nodes may only have a few neighbors. Intuitively, the network size  $n$  should not be small and most of communications between nodes and the base station are multi-hop; otherwise, there is no point for introduction of data aggregation. Our protocol does not rely on specific data aggregation modes; any effective aggregation organization, no matter it is cluster-, chain-, tree-, or grid-based [3], suffices for the proposed protocol. We do not address the problem how to construct a specific aggregation organization, because it is quite irrelevant to our proposal.

There are two kinds of roles for sensor nodes: a *contributor* that detects events and generates *raw messages*, and an *aggregator* that aggregates all messages that it received plus possibly its own raw messages<sup>1</sup> and then forwards *aggregated messages*. Messages are aggregated though the network; and the base station eventually retrieves a number of aggregated messages, which are defined as *reports*.

We assume that there is a secure broadcast message authentication scheme available for the base station. We may use the popular, lightweight, symmetric-key hash-chain-based  $\mu$ TESLA authentication protocol [22] for this purpose. Alternatively, all nodes are preloaded with the base station's public key, and then messages from the base station are signed and verified through classical signature schemes.

Every node shares a distinct long-term key with the base station. Those keys are pre-distributed into sensor nodes prior to deployment. No pairwise keys between sensor nodes or any other kinds of keys are demanded. Sensor nodes are equipped with a cryptographic hash function. Those constitute all security primitive requirements of the proposed protocol. As a general method to reduce its key storage, the base station may employ a pseudo-random generator (PRG) along with one main key, which generates node's key using a node's ID and the main key as input of PRG.

### C. Security Objectives

Data confidentiality is not a concern for the application scenarios which our proposal is targeted at. For simple event detection networks, if an adversary is physically close to sensor nodes, he can easily sense those simple events on his own. Certainly, some extent of message privacy is appreciated. Barely given a message and without any other knowledge, should an adversary not tell which nodes have contributed to the message and what events are reported.

Generally speaking, the main security goal of secure aggregation is to strictly restrain adversary's influence only on those compromised nodes. First, even with collusion of all other sensor nodes, an adversary cannot fake a specific node's event report with more than a low, theoretically-limited probability. Second, other than compromised nodes' inputs on event report, the adversary is only able to affect accumulative results in a trivial way. Otherwise, the scheme should detect the abnormality and discard related results. If intermediate nodes manipulate or drop integrity nodes' event-report messages, the secure aggregation scheme should be able to effectively detect the attacks.

### D. Native Solution without Data Aggregation

There exists a straightforward, native solution to fulfilling the aforementioned security objectives, which does not engage data aggregation and then incurs heavy communication cost. When node  $u$  detects an event  $e$ , it generates a raw message composed of  $\text{Mac}(K_u, e)$  and sends to the base station, where  $K_u$  is the symmetric key shared by node  $u$  and the base station, and  $\text{Mac}$  denotes a tag by a standard deterministic MAC. When the base station receives a message, it may try all combinations to figure out which node claims which event. In order to prevent message-dropping by malicious nodes, the base station broadcasts all messages it received through the broadcast authentication mechanism. If a node does not obtain the authenticated acknowledgment of its raw message in a reasonable time slot, it will notify the base station about a potential of the attack.

It may seem tempting to use aggregate MAC [23], which aggregates multiple MAC tags by different contributors into a single tag that can be verified by the base station, to reduce the communication payload of the native scheme. Unfortunately, in such a application, an aggregate tag has to list all contributors and associated events such that the base station can verify the tag. Consequently, the size of an aggregated message is still proportional to the number of contributors and events. Moreover, it violates the message privacy objective.

## IV. SECURE BLOOM FILTER

### A. Bloom Filter

Bloom filter, conceived by Bloom [6], is a simple, space-efficient probabilistic data structure that succinctly represents a set in order to support membership queries. Due to its distinguished space advantages and excellent distributed properties, bloom

<sup>1</sup>In such a case, the aggregator is also a contributor.

filter has been widely used in numerous areas, such as web cache sharing [24], distributed storage system [25].

Typically, a Bloom filter is implemented as a bit-array of  $m$  bits associated with  $h$  different hash functions, each of which maps a set element to one of the  $m$  array positions in a uniformly random manner. All bits in an initial Bloom filter are set to 0, standing for an empty set. To insert an element  $u$  into a set represented by a Bloom filter BF,  $h$  array positions are calculated by hash functions on  $u$  and the bits at those positions in BF are set to 1. Correspondingly, when it is required to check the membership of an element  $v$  within the Bloom filter BF, supplying  $v$  to hash functions outputs  $h$  array positions; if any of the bits at the  $h$  positions is 0, then element  $v$  does not belong to the set; otherwise, the element is claimed to be a member of the set.

It is clear that there is no *false negative* in the Bloom filter membership verification—an element which tests negative within a Bloom filter definitely is not a legitimate member of the set. On the other hand, Bloom filter may yield *false positive*: a member outside the set passing the membership verification on the Bloom filter. The probability of a false positive for an element not in the set, or the *false positive probability*, can be calculated in a straightforward manner.

Let  $\zeta$  be the probability of a bit being 0 in a Bloom filter, the false positive probability is then

$$P_{\text{FP}} = (1 - \zeta)^h . \quad (1)$$

After inserting  $t$  elements in a Bloom filter, assume that the hashing functions outputs are independently uniformly distributed, we have

$$\zeta = \left(1 - \frac{1}{m}\right)^{ht} \approx e^{-ht/m} ;$$

the false positive probability is therefore

$$P_{\text{FP}} = \left(1 - \left(1 - \frac{1}{m}\right)^{ht}\right)^h \approx \left(1 - e^{-ht/m}\right)^h . \quad (2)$$

The right hand side is minimized when

$$h = \frac{m}{t} \ln 2 \approx 0.6931 \frac{m}{t} .$$

In such a case,  $\zeta = 0.5$ , and the false positive probability is  $2^{-h} \approx (0.6185)^{m/t}$ . Hence, the optimal results are achieved when each bit of the Bloom filter is 0 with probability 1/2.

In a typical parameters determination procedure of Bloom filter, users specify desired false positive probability  $P_{\text{FP}}$  and predicted inserted element number  $t$ , then  $h$  and  $m$  can be calculated via

$$h = \lceil -\log_2 P_{\text{FP}} \rceil$$

and

$$m = \lceil ht / \ln 2 \rceil .$$

### B. Technical Core: Secure Bloom Filter

The technical core of our proposal is a variant of Bloom filter, for which we coin the term *secure Bloom filter*, by deploying a cryptographic hash function to substitute a family of hash functions for Bloom filter. We use the cryptographic hash function by means of HMAC because it is provably secure [26] and has been standardized [27]. Let  $H(*)$  denote a cryptographic hash function; according to RFC 2104 [27], to authenticate a message  $m$  with a secret symmetric key  $K$ ,  $\text{HMac}(K, m)$  is defined by

$$\text{HMac}(K, m) = H((K \oplus \text{opad}) \parallel H((K \oplus \text{ipad}) \parallel m)),$$

where  $\parallel$  denotes concatenation,  $\oplus$  denotes exclusive-OR (XOR), opad is the one-block-size outer padding (0x5c5c...5c), and ipad is the one-block-size inner padding (0x3636...36).

Now we make use of HMAC to construct secure Bloom filter that efficiently supports event reporting and message aggregation in wireless sensor networks. Suppose there are  $c$  pre-defined events  $(e_1, e_2, \dots, e_c)$ , and each event is associated with an event identifier, which is a nonce (number used only once) and may be refreshed by the base station on a regular basis. When node  $u$  with identity  $\text{id}_u$  detects an event  $e_i$ , whose current event identifier is  $N_{e_i}$ , the node creates a secure Bloom filter  $\text{BF}_{u, e_i}$  as follows.

For simplicity and computational efficiency, we assume that the size of a secure Bloom filter is a power of two, say  $m = 2^d$ , where  $d$  is an integer. We stress that in principal  $m$  can be any positive integer. Recall that node  $u$  shares a distinct symmetric key  $K_u$  with the base station. The node calculates  $\text{HMac}(K_u, \text{id}_u \parallel N_{e_i})$  and the output is divided into  $k$  equal size pieces of  $d$  bits each <sup>2</sup> such that each piece indicates a position on the Bloom filter array. We denote this processing

<sup>2</sup>If the output block size of HMAC is greater than  $dh$ , it can be truncated. When the size is less than  $dh$ , we may conduct several rounds of HMAC operations with different salts to obtain sufficient output.

TABLE I: Secure Bloom filter parameter lookup table

$h$	7	8	9	10	11	12	13	14	15	16
$i$ for maximizing $\tau(h, i)$	5	6	7	8	8	9	10	10	11	12
maximal $\tau(h, i)$	16800	191520	2.33E+6	3.02E+7	4.79E+08	8.08E+09	1.43E+11	2.73E+12	5.91E+13	1.32E+15
min( $m$ ) for $\text{Adv} \leq 2^{-80}$	11066	4684	2417	1434	952	680	514	406	334	282
min( $m$ ) for $\text{Adv} \leq 2^{-100}$	80176	26495	11278	5733	3356	2159	1492	1093	842	671

by  $\text{BF}_{u,e_i} \leftarrow \text{HMac}(K_u, id_u || N_{e_i})$ . Based on the pseudo-randomness assumption of the cryptographic hash function  $H(*)$ , it effectively performs as  $k$  distinct hash functions with uniformly-distributed output, and all theoretical analysis on the classical Bloom filter still holds for the secure Bloom filter.

*Definition 1:* The advantage of an adversary against secure Bloom filter is defined as the probability that the adversary successfully constructs a valid secure Bloom filter associated with a node claiming an event, without the knowledge of the node's key. Let  $\text{Adv}$  denote the advantage.

*Derivation of Adv:* Based on the security of HMAC and the pseudo-randomness of cryptographic hash function, we derive the advantage of an adversary against secure Bloom filter with parameters  $(m, h)$ . Intuitively, for 80-bit security, the selections of  $(m, h)$  should satisfy  $\text{Adv}(m, h) \leq 2^{-80}$ . Without loss of generality, we suppose that adversary  $\mathcal{A}$  is required to output a secure Bloom filter  $\text{BF}^{\mathcal{A}}$  which should be identical to  $\text{BF}_{u,e}$  for node  $u$  claiming event  $N_e$ .

Let  $p_1(i)$  be the probability of the number of bit one in  $\text{BF}_{u,e}$  being  $i$ . Obviously  $p_1(i)$  is none-zero only if  $1 \leq i \leq h$ . Based on the pseudo-randomness of cryptographic hash function, we have

$$p_1(i) = \frac{\binom{m}{i} \tau(h, i)}{m^h},$$

where  $\tau(h, i)$  is equal to the output of the following problem: How many distinct words we can get if we use  $i$  different letters to construct words of length  $h$  while it is required that all  $i$  letters are used.

If HMAC is secure, then the adversary only can output random  $\text{BF}^{\mathcal{A}}$  with the number of bit one that he selects. When the bit one number of  $\text{BF}_{u,e}$  is  $i$ , the probability that a random  $\text{BF}^{\mathcal{A}}$  with  $i$  bits of one equals  $\text{BF}_{u,e}$  is then

$$p_2(i) = \frac{1}{\binom{m}{i}}.$$

Subsequently, the overall advantage that adversary can achieve is

$$\text{Adv}(m, h) = \max_{1 \leq i \leq h} p_1(i) p_2(i) = m^{-h} \max_{1 \leq i \leq h} \tau(h, i). \quad (3)$$

Consider  $\tau(h, i)$  in the context of the word-counting problem, it is easy to get the recursive formula:

$$\tau(h, i) = \begin{cases} i^h - \sum_{j=1}^{i-1} \binom{i}{j} \tau(h, j) & \text{if } i > 1 \\ 1 & \text{if } i = 1 \end{cases}. \quad (4)$$

In practice,  $h$  is relatively small, thus we just explore all values of  $\tau(h, i)$  for  $1 \leq i \leq h$  by Eq. (4), and then retrieve  $\max_{1 \leq i \leq h} \tau(h, i)$ . Subsequently,  $\text{Adv}(m, h)$  can be computed by Eq. (3), or we can determine the minimal  $m$  satisfying a designated security level. To facilitate parameter selections, we calculate and form Table I as secure Bloom filter parameters lookup table for  $h = 7, 8, \dots, 16$ .

It is worth noting that, by applying Stirling's approximation on

$$p_2(h) = \frac{1}{\binom{m}{h}} = \frac{(m-h)! h!}{m!},$$

we can get

$$p_2(h) \approx h! \left(\frac{e}{m}\right)^h \left(1 - \frac{h}{m}\right)^{m-h+0.5},$$

which may be further approximated as

$$p_2(h) \approx h! e^{(h^2-0.5h)/m} m^{-h} \approx h! m^{-h},$$

when  $m \gg h$ . In other words, even if an adversary is aware that a particular secure Bloom filter has  $h$  bits of ones, the probability that he successfully forges it is roughly  $h! m^{-h}$ .

For the specific choice of the cryptographic hash function, at present we select SHA-1. Even though collisions have been found in SHA-1 [28], which raised serious concerns on its security, virtually it does not damage the secure Bloom filter, since

---

**Algorithm 1** *ProcessReport*(BF<sub>j</sub>)

---

```

1: if the bit one percentage of BFj exceeds  $\rho$  then
2:   ALERT: BFj is corrupted, discard it
3:   return false
4: end if
5: create an empty BFtest for testing
6: for each BFu,ei do
7:   if BFu,ei  $\wedge$  BFj = BFu,ei then
8:     determines that node  $u$  has claimed event  $e_i$ 
9:     BFtest  $\leftarrow$  BFtest  $\vee$  BFu,ei
10:  end if
11: end for
12: if BFtest  $\neq$  BFj then
13:   ALERT: BFj is corrupted, discard it
14:   cancel off all previous determinations
15:   return false
16: end if
17: return true

```

---

its security, similar to HMAC [26], relies on second preimage resistance of underlying hash function, rather than the more strong requirement of collision resistance. By now, there is no published result of second preimage attacks against SHA-1. Of course, when SHA-3 is finalized in the future, SHA-1 may be replaced by SHA-3 in the protocol for the purpose of higher security level.

## V. PROPOSED PROTOCOL

### A. Protocol Description and Analysis

An interesting property of Bloom filter is inherently supporting aggregation by bitwise-ORing Bloom filters of a same kind together to generate one Bloom filter such that the set represented by the output Bloom filter is the union of the sets of input Bloom filters. Let  $\vee$  denote the bitwise-OR operation on Bloom filters. This desired aggregation property perfectly applies to secure Bloom filters generated by sensor nodes with different keys, and thus we take advantage of it in the proposed protocol, which consists of the following four phases.

1) *Initialization*: To activate sensor nodes detection on selected events  $(e_1, e_2, \dots, e_c)$ , the base station broadcasts to all nodes an authenticated task message, including fresh nonces  $(N_{e_1}, N_{e_2}, \dots, N_{e_c})$  as event identifies and events specifications if necessary.

2) *Event Claiming and Aggregation*: When node  $u$  detects that event  $N_{e_i}$  occurs, it generates a new Bloom filter  $\text{BF}_{u,e_i} \leftarrow \text{HMac}(K_u, id_u \| N_{e_i})$ , which constitutes the raw message by node  $u$  as a contributor to claim event  $N_{e_i}$ .

The aggregation on messages is straightforward and efficient. Suppose an aggregator node  $v$  receives  $j$  messages of  $\text{BF}_i^v$ , where  $i = 1, 2, \dots, j$ , some of which might be contributed by node  $v$  itself, it calculates an aggregated message by

$$\text{BF}' = \text{BF}_1^v \vee \text{BF}_2^v \vee \dots \vee \text{BF}_j^v,$$

To maintain a reasonably slow false positive probability of Bloom filter, aggregators cannot aggregate messages without a limitation. As we see, the way that we use Bloom filter is quite different from most of other applications. As the message carrier to traverse networks, Bloom filter should not use large  $m$  for the sake of communication cost. In fact,  $m \leq nc$  for most cases. As result, in addition to  $m$  and  $h$  for Bloom filter settings, our protocol requires a third parameter: bit one percentage upper-bound  $\rho$ . Accordingly, aggregator  $v$  tries its best to output one or several aggregated messages with respect to the threshold  $\rho$ . In that way, there would be no legitimate Bloom filter whose bit one percentage exceeds  $\rho$ . Recall that the false positive probability is optimized when the probability of a bit being 0 in a Bloom filter is 1/2 (see Section IV-A), the threshold  $\rho$  is set to be 50% for most circumstances.

After aggregating messages with respect to  $\rho$ , aggregator  $v$  sends the aggregated message(s) to next node according to the underlying data aggregation organization.

3) *Processing Reports*: Eventually, the base station receives  $l$  reports  $\text{BF}_j$ , where  $j = 1, 2, \dots, l$ . Using its distinct keys shared with nodes, the base station calculates all  $nc$  Bloom filters corresponding to node-event pairs beforehand. For each report, the base station conducts Algorithm 1 to proceed it, intended to discover which nodes claimed which events and to detect attacks.

Let  $\bigwedge$  denote bitwise-AND operation, during the processing procedure, if

$$BF_{u,e_i} \bigwedge BF_j = BF_{u,e_i} ,$$

that is, all bit one occurrences in  $BF_{u,e_i}$  appear in  $BF_j$ , the base station determines that node  $u$  has claimed event  $e_i$ . Certainly, false positives may occur for this verification as an inherent downside of Bloom filter. We define the *false positive rate*  $R_{FP}$  as the ratio of the count of false positives over the node-event pair number (i.e.  $nc$ ).

In the ideal case where all  $l$  reports are independent and have  $\rho$  proportion of bit 1, the false positive probability within a report is  $\rho^h$  by Eq. (1). Suppose  $nc$  is much greater than the raw message number, then the number of false positives roughly follows a binomial distribution of parameter  $(nc, 1 - (1 - \rho^h)^l)$ , and then its expected value approximates  $nc(1 - (1 - \rho^h)^l)$ . Therefore,

$$R_{FP}^{(\text{ideal})} \approx 1 - (1 - \rho^h)^l \approx 1 - e^{-l\rho^h} . \quad (5)$$

In reality, most of reports have less than  $\rho$  fraction of bit one, and the expected value of false positive number is related to  $l$  as bigger  $l$  may indicate more number of raw messages. Therefore,  $R_{FP}^{(\text{ideal})}$  is an upper bound of  $R_{FP}$ .

If compromised nodes misleadingly claim events with their own keys, clearly there are no effective methods to prevent this attack except for designing and conducting complicated and expensive detection protocols. We are not concerned about that issue; instead the proposed protocol is supported to detect message manipulation attack. The steps related to  $BF_{\text{test}}$  in Algorithm 1 carry out that mission.

Recall that by selecting proper  $(m, h)$ , it is infeasible for the adversary to construct a legitimate secure Bloom filter associated with integrity nodes. During transmission, messages may be maliciously altered by adding or deleting some bits. To detect this kind of attack, actually we only need to consider the case that some bits in reports are *randomly* flipped by the adversary. We stress that the base station is solely responsible and capable of processing messages and those effective manipulations shall eventually reflect on reports; if the adversary replaces a message by another message, in fact it is an interception rather than manipulation, and the countermeasure will be addressed later. If the adversary adds a valid message to an unrelated report, that is not a manipulation attack either, and merely it slightly increase the false positive rate, which is still bounded by Eq. (5).

We evaluate the probability  $P_{\text{det}}$  of the base station detecting that a random bit zero in a report is flipped to bit one, which might be the manipulation case that is detected with least probability. Using the previous ideal case of studying  $R_{FP}$ , we can get

$$\begin{aligned} P_{\text{det}}^{(\text{ideal})} &\approx \left( 1 - \rho^h \left( 1 - \left( 1 - \frac{1}{\rho m} \right)^h \right) \right)^{nc} \\ &\approx \left( 1 - \rho^h + \rho^h e^{-h/\rho m} \right)^{nc} . \end{aligned} \quad (6)$$

4) *Subsequent Verification by All Nodes:* To detect if malicious nodes intercepted messages related to intact nodes, the base station broadcasts those  $l$  reports throughout the whole network. In addition, in order to prevent the message replay attack, the base station should notify all sensor nodes to update event identifiers with new nonces. Surely this stage can be twisted with the initialization stage to add/update/delete events. All those messages are still authenticated by the base station's broadcast authentication scheme.

As long as an intact node is not isolated by compromised nodes from the base station, the node can receive those reports and verify the authenticity of messages. Therefore, the node that claimed the event can determine whether its report is included in the final results. If not, the node notifies this abnormality to the base station in any available secure, uni-cast means by its node key. In addition, if nodes find that there are some false positives related to them, they may report to the base station. If the base station determines that the number of missing raw messages exceeds a reasonable threshold, it may conclude that there exists a message interception attack in the sensor networks and further security mechanisms can be conducted.

## B. Compression and Parameters

1) *Compression on Bloom Filter:* The entropies of raw messages are pretty low—most of bits are zero because there is only one element in a set that the raw message Bloom filter represents. As a natural way to reduce the transmission size of Bloom filter, this kind of Bloom filters should be compressed in our protocol. The problem of compressed Bloom filter has been researched by Mitzenmacher [29], who demonstrates that a better compression rate can be achieved by utilizing smaller  $h$  and bigger  $m$  while maintaining the same level of false positive probability.

However, that rationale does not fit the requirement of our protocol. First, the selection of  $h$  matters to the security level of secure Bloom filter. Second, there are many Bloom filters with varying entropies during the protocol procedure, one pre-determined parameter set hardly can provide overall desired transmission savings. Therefore, we does not make use of the results in [29]. Instead, a succinct compression mechanism is employed in our protocol: one additional bit  $b$  attached to a Bloom filter indicates that the filter is normally stored as before when  $b = 0$ , and if  $b = 1$ , then the filter is expressed as a

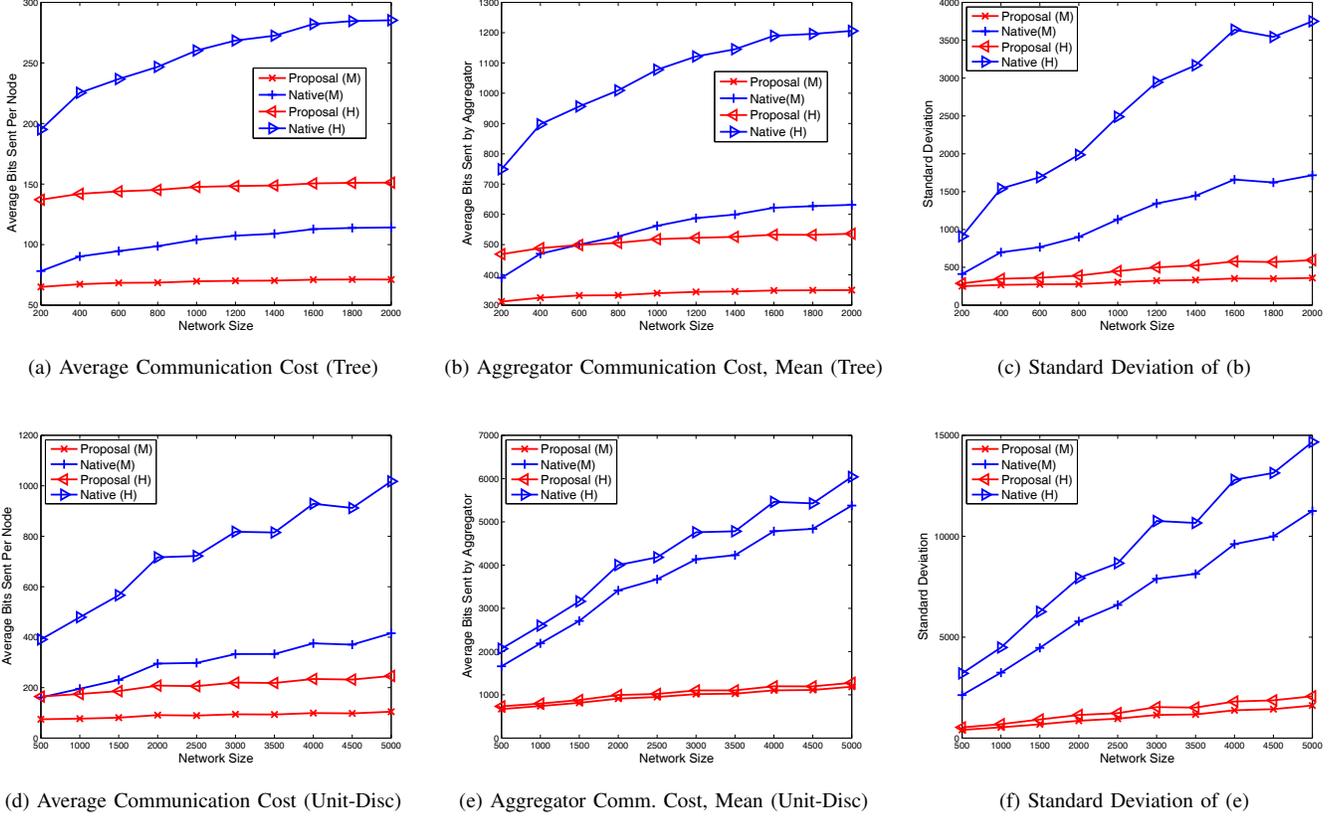


Fig. 1: Simulation Results on Protocols Communication Performances

list of all bit one positions. Since the number of bit one is dynamical, such a list should contain a number-indication segment of bit length  $w = \lceil \log_2(m/d) \rceil$ , where  $m = 2^d$ . Before transmitting a Bloom filter, a node chooses the coding which leads to smaller transmission size. Specifically, the list expression is used only if the number of bit one is less than  $\frac{m-w}{d}$ .

2) *Parameters, Computation and Storage Costs*: For sensor networks of size  $n$  varying from hundreds to thousands, and to provide 80-bit security level, we recommend  $m = 1024$ ,  $h = 11$  according to Table I. Because we use SHA-1, whose output size is 160 bit, performing one HMAC function that consists of two SHA-1 calculations would suffice for nodes to construct one raw message. We assume that the key length is 160-bit, so the primary storage cost of our protocol is merely 20 bytes. Therefore, our protocol is computation- and storage-efficient.

## VI. SIMULATIONS

To verify our theoretical analysis and evaluate the performances, we implement the proposed data aggregation protocol and run simulations on the OMNeT++ framework [30].

### A. Simulation Design

1) *Network Scenarios*: We run simulations in two network scenarios. The first one is an abstract hierarchical network modeled as a *Tree* rooted on the base station, which fairly emulates many wireless sensor networks. We generate a tree topology in a way that the number of a non-leaf node's children is uniformly selected from four, five, and six. In our experiments, the network size for the Tree Topology increases from 200 to 2000 at a step of 200.

The other one is *Unit-Disc Graph*, in which nodes are uniformly deployed in a  $500 \times 500$  square and nodes follow the standard unit-disc bidirectional communication model. We run simulations under the Unit-Disc Graph for network sizes ranging from 500 to 5000 with a space of 500 between. In order to imitate a dense sensor network and maintain consistence, for networks of different sizes, we adjust the node communication range such that the average node degree keeps the approximate 20. This kind of wireless sensor network simulation scenario models random deployed sensor networks and has been widely used in many literatures, such as [31].

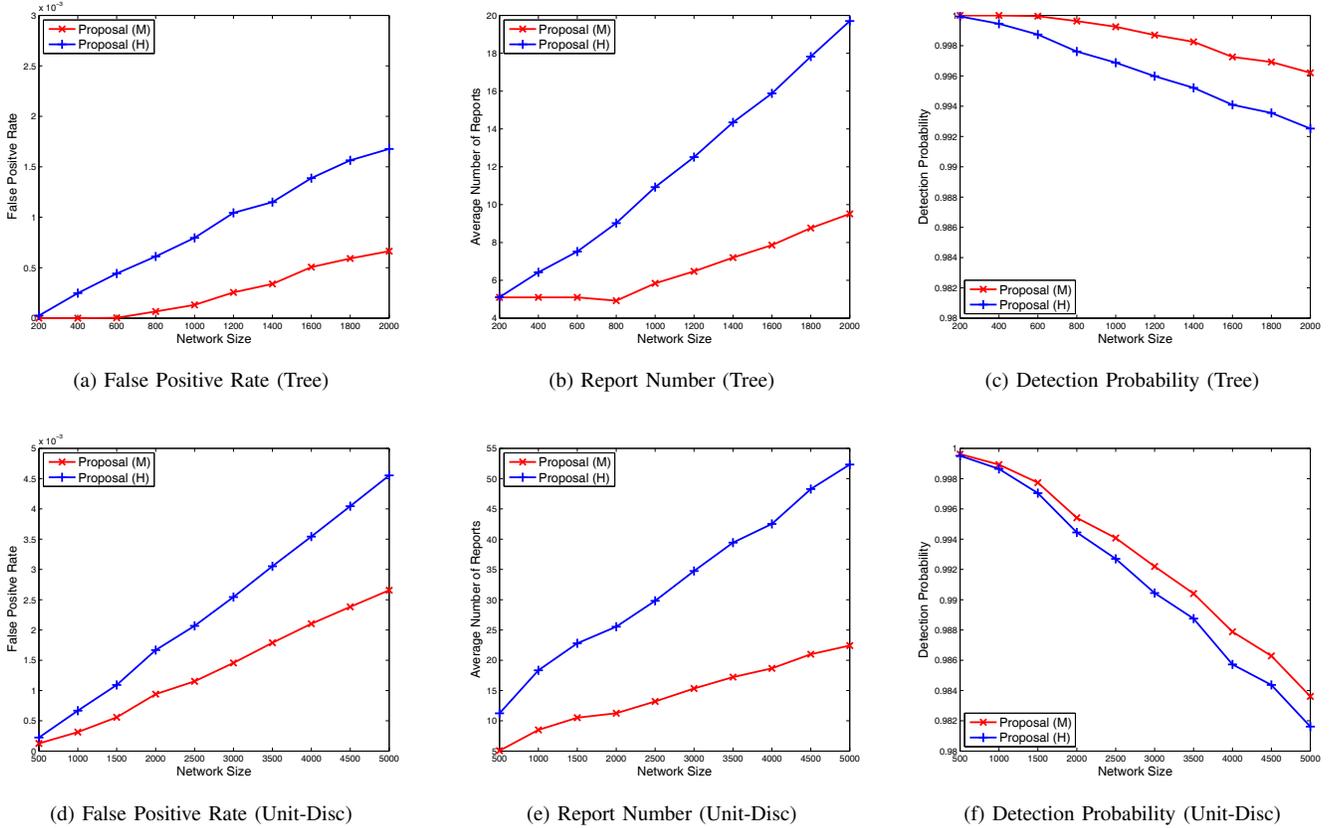


Fig. 2: Simulation Results Related to Protocol Security

By standard software engineering methodologies, the implemented layer of the data aggregation protocol is encapsulated and separated from underlying network topologies. With a little effort on topology implement, irregular surveillance areas and other network scenarios, such as random graph and random tree, can be seamlessly introduced to run simulations.

2) *Mechanisms of Claiming Events*: On each network scenario, we specify two different workloads to activate nodes claiming events, one for mild work burden, and the other for relatively heavy turnover. For the Tree Topology, in each round of protocol procedure, 20% and 50% of sensor nodes, respectively for two levels of workload, are randomly chosen as contributors to claim an event. For the Unit-Disc Graph, we designate that nodes are capable of detecting events taking place no farther than a distance of 100, and at the beginning of each round, an event occurs at two and six positions (respectively for two workloads and analogue to 20% and 50% contributors settings in the Tree topology) that are uniformly selected in the square.

3) *Simulation Arrangement*: Overall, there are four test cases: Tree Topology with mild workload, Tree Topology with heavy workload, Unit-Disc Graph with mild workload, and Unit-Disc Graph with heavy workload. Each case is associated with ten different network sizes, requiring 40 experiments. In order to obtain relatively fair and comparable results, ten random networks in accordance with the parameter settings are constructed for each experiment, each of whose simulation executions is quoted as a run; and in each run, we perform 100 rounds of protocol processing with new events and randomly selected contributors/event-locations. For the purpose of communication cost comparison, in addition to our proposed protocol, the native scheme described in Section III-D is also conducted under the identical conditions. We assume that the message length is 128-bit in the native scheme. Since the last stage is needed for both our protocol and the native scheme, the experiment results on communication costs do not consider that stage. To evaluate the protocol's detection probability  $P_{det}$ , we simulate the manipulation attack by randomly flipping a bit zero of a report and then measure whether the base station subsequently detects that the report is illegitimate. For each report, this procedure is repeated by twenty times. The final runs outputs are averaged to create the corresponding experiment results.

## B. Experimental Results

In figures of this paper, "Proposal", "Native", "(M)", and "(H)" respectively stand for our proposed data aggregation protocol, the native scheme, mild workload, and heavy workload. The protocols communication performances along with their comparison

are depicted in Figure 1. For general sensor network protocols, the communication cost is often represented by average transmission bits per node (Fig. 1a and Fig. 1d). For a data aggregation scheme, the energy consumption of aggregators and the balance among them are critical performance metrics; so we also measure the means (Fig. 1b and Fig. 1e) and the standard deviations (Fig. 1c and Fig. 1f) of bits sent by aggregators (intermediate nodes for the native scheme). The simulation results clearly demonstrate that for all four test cases, our protocol remarkably outweighs the native scheme on all three metrics; not only the protocol reduces the overall energy consumption and reserves previous energy of aggregators, but also it achieves splendid balance among aggregators, which is highly desired for data aggregation scheme. Moreover, as those measurements for our protocol gently grow with  $n$ , it proves that the proposed protocol scales to network size well, which is quite appreciated for large-size sensor networks.

The simulation results regarding the protocol's security are illuminated in Fig. 2. Specifically, Fig. 2a and Fig. 2d depict the protocol false positive rates; Fig. 2b and Fig. 2e show the average report numbers; Fig. 2c and Fig. 2f exhibit the detection probabilities. The false positive rates of our protocol are fairly small, compared to typical Bloom filter applications. By applying the experimental results of report number  $l$  to Eq. (5), we verify that the previous theoretical analysis on  $R_{FP}$  holds as all false positive rates are indeed bounded on Eq. (5). Lastly, the high detection probabilities, all of which are greater than 98% in the simulations, exemplify the strong security of our protocol.

## VII. CONCLUSION

In this paper, we propose a practical secure data aggregation protocol for wireless sensor network based on Bloom filter, and conduct thorough theoretical analysis on the related topics. The extensive simulation results sufficiently demonstrate that the proposal presents remarkable performances on communication cost, energy consumption balance, and security.

## ACKNOWLEDGMENT

The research is supported by NSERC Strategic Project Grants.

## REFERENCES

- [1] D. Estrin, R. Govindan, J. Heidemann, and S. Kumar, "Next Century Challenges: Scalable Coordination in Sensor Networks," in *Proceedings of the 5th ACM/IEEE International Conference on Mobile Computing and Networking*. Seattle: IEEE Computer Society, 1999, pp. 263–270.
- [2] I. F. Akyildiz, W. Su, Y. Sankarasubramaniam, and E. Cayirci, "A survey on sensor networks," *IEEE Communications Magazine*, vol. 40, no. 8, pp. 102–114, 2002.
- [3] R. Rajagopalan and P. K. Varshney, "Data-aggregation techniques in sensor networks: a survey," *IEEE Communications Surveys & Tutorials*, vol. 8, no. 4, pp. 48–63, 2006.
- [4] E. Fasolo, M. Rossi, J. Widmer, and M. Zorzi, "In-network aggregation techniques for wireless sensor networks: a survey," *IEEE Wireless Communications*, vol. 14, no. 2, pp. 70–87, 2007.
- [5] P. Baronti, P. Pillai, V. W. C. Chook, S. Chessa, A. Gotta, and Y. F. Hu, "Wireless sensor networks: A survey on the state of the art and the 802.15.4 and ZigBee standards," *Computer Communications*, vol. 30, no. 7, pp. 1655–1695, 2007.
- [6] B. H. Bloom, "Space/time trade-offs in hash coding with allowable errors," *Communications of the ACM*, vol. 13, no. 7, pp. 422–426, 1970.
- [7] D. Westhoff, J. Girao, and M. Acharya, "Concealed Data Aggregation for Reverse Multicast Traffic in Sensor Networks: Encryption, Key Distribution, and Routing Adaptation," *IEEE Transactions on Mobile Computing*, vol. 5, no. 10, pp. 1417–1431, 2006.
- [8] C. Castelluccia, A. C. F. Chan, E. Mykletun, and G. Tsudik, "Efficient and provably secure aggregation of encrypted data in wireless sensor networks," *ACM Trans. Sen. Netw.*, vol. 5, no. 3, pp. 1–36, 2009.
- [9] P. Steffen, W. Dirk, and C. Claude, "A Survey on the Encryption of Convergecast Traffic with In-Network Processing," *IEEE Transactions on Dependable and Secure Computing*, vol. 7, pp. 20–34, 2010.
- [10] C. Fontaine and F. Galand, "A survey of homomorphic encryption for nonspecialists," *EURASIP Journal on Information Security*, vol. 2007, no. 1, pp. 1–10, 2007.
- [11] C. Gentry, "Fully homomorphic encryption using ideal lattices," in *Proceedings of the 41st annual ACM symposium on Theory of computing*. Bethesda, MD, USA: ACM, 2009, pp. 169–178.
- [12] M. van Dijk, C. Gentry, S. Halevi, and V. Vaikuntanathan, "Fully Homomorphic Encryption over the Integers," in *EUROCRYPT'10*, 2010.
- [13] Z. Li and G. Gong, "Data Aggregation Integrity Based on Homomorphic Primitives in Sensor Networks," in *Proceedings of 9th International Conference on Ad Hoc Networks and Wireless (ADHOC-NOW 2010)*. LNCS 6288, 2010, pp. 149–162.
- [14] L. Hu and D. Evans, "Secure aggregation for wireless networks," in *Proceedings of the 2003 Symposium on Applications and the Internet Workshops (SAINT'03 Workshops)*, 2003, pp. 384 – 391.
- [15] P. Jadia and A. Mathuria, "Efficient Secure Aggregation in Sensor Networks," in *High Performance Computing (HiPC 2004)*. LNCS 3296, 2004, pp. 40–49.
- [16] Y. Yang, X. Wang, S. Zhu, and G. Cao, "A Secure Hop-by-Hop Data Aggregation Protocol for Sensor Networks," in *Proceedings of the 7th ACM international symposium on Mobile ad hoc networking and computing (MobiHoc'06)*, 2006, pp. 356 – 367.
- [17] B. Przydatek, D. Song, and A. Perrig, "SIA: Secure Information Aggregation in Sensor Networks," in *Proceedings of the first international conference on Embedded networked sensor systems (SenSys'03)*, Los Angeles, California, USA, 2003, pp. 255 – 265.
- [18] H. Chan, A. Perrig, and D. Song, "Secure hierarchical in-network aggregation in sensor networks," in *Proceedings of the 13th ACM conference on Computer and communications security*. Alexandria, Virginia, USA: ACM, 2006, pp. 278–287.
- [19] K. B. Frikken and J. A. Dougherty I. V., "An efficient integrity-preserving scheme for hierarchical sensor aggregation," in *Proceedings of the first ACM conference on Wireless network security (WiSec'08)*. Alexandria, VA, USA: ACM, 2008, pp. 68–76.
- [20] H. Chan and A. Perrig, "Efficient security primitives derived from a secure aggregation algorithm," in *Proceedings of the 15th ACM conference on Computer and communications security (CCS'08)*. Alexandria, Virginia, USA: ACM, 2008.
- [21] Y. Liu, C. Li, Y. He, J. Wu, and Z. Xiong, "A Perimeter Intrusion-Detection System Using Dual-Mode Wireless Sensor Networks," in *Second International Conference on Communications and Networking in China (CHINACOM '07)*, 2007, pp. 861–865.
- [22] A. Perrig, R. Szewczyk, V. W. D. Culler, and J. D. Tygar, "SPINS: Security protocols for sensor networks," in *Proceedings of the Annual International Conference on Mobile Computing and Networking (MOBICOM)*. Rome Italy: IEEE, 2001, pp. 189–199.

- [23] J. Katz and A. Lindell, "Aggregate Message Authentication Codes," in *CT-RSA 2008*. LNCS 4964, 2008, pp. 155–169.
- [24] L. Fan, P. Cao, J. Almeida, and A. Z. Broder, "Summary cache: a scalable wide-area web cache sharing protocol," *IEEE/ACM Transactions on Networking*, vol. 8, no. 3, pp. 281–293, 2000.
- [25] F. Chang, J. Dean, S. Ghemawat, W. C. Hsieh, D. A. Wallach, M. Burrows, T. Chandra, A. Fikes, and R. E. Gruber, "Bigtable: a distributed storage system for structured data," in *Proceedings of the 7th USENIX Symposium on Operating Systems Design and Implementation*, vol. 7. Seattle, WA: USENIX Association, 2006, pp. 15–15.
- [26] M. Bellare, R. Canetti, and H. Krawczyk, "Keying Hash Functions for Message Authentication," in *Advances in Cryptology - CRYPTO'96*. LNCS 1109, 1996, pp. 1–15.
- [27] H. Krawczyk, M. Bellare, and R. Canetti, "RFC 2104 - HMAC: Keyed-Hashing for Message Authentication," IETF, 1997.
- [28] X. Wang, Y. L. Yin, and H. Yu, "Finding Collisions in the Full SHA-1," in *Advances in Cryptology - CRYPTO 2005*. LNCS 3621, 2005, pp. 17–36.
- [29] M. Mitzenmacher, "Compressed bloom filters," *IEEE/ACM Transactions on Networking*, vol. 10, no. 5, pp. 604–612, 2002.
- [30] A. Varga and R. Hornig, "An overview of the OMNeT++ simulation environment," in *Proceedings of the 1st international conference on Simulation tools and techniques for communications, networks and systems & workshops*. Marseille, France: ICST, 2008, pp. 1–10.
- [31] B. Parno, A. Perrig, and V. Gligor, "Distributed Detection of Node Replication Attacks in Sensor Networks," in *Proceedings of the IEEE Symposium on Security and Privacy*, 2005, pp. 49 – 63.