

# Secure Aggregation in a Publish-Subscribe System

Kazuhiro Minami<sup>†</sup>, Adam J. Lee<sup>‡</sup>, Marianne Winslett<sup>†</sup>, and Nikita Borisov<sup>†</sup>  
minami@cs.uiuc.edu, adamlee@cs.pitt.edu, winslett@cs.uiuc.edu, nikita@uiuc.edu

<sup>†</sup>University of Illinois at Urbana-Champaign

<sup>‡</sup>University of Pittsburgh

## ABSTRACT

A publish-subscribe system is an information dissemination infrastructure that supports many-to-many communications among publishers and subscribers. In many publish-subscribe systems, in-network aggregation of input data is considered to be an important service that reduces the bandwidth requirements of the system significantly. In this paper, we present a scheme for securing the aggregation of inputs to such a publish-subscribe system. Our scheme—which focuses on the additive aggregate function *sum*—preserves the confidentiality and integrity of aggregated data in the presence of untrusted routing nodes. Our scheme allows a group of publishers to publish aggregate data to authorized subscribers without revealing their individual private inputs to either the routing nodes or the subscribers. In addition, our scheme allows subscribers to verify that routing nodes perform the aggregation operation correctly. We use a message authentication code (MAC) scheme based on the discrete logarithm property to allow subscribers to verify the correctness of aggregated data without receiving the digitally-signed raw data used as input to the aggregation. In addition to describing our secure aggregation scheme, we provide formal proofs of its soundness and safety.

**Categories and Subject Descriptors:** C.2.4 [Distributed Systems]: Distributed applications; K.6.5 [Management of Computing and Information Systems]: Security and Protection

**General Terms:** Security

**Keywords:** Aggregation, data privacy, integrity, publish-subscribe systems

## 1. INTRODUCTION

A publish-subscribe (pub-sub) system [2, 3, 21, 25] is an information dissemination infrastructure that supports many-to-many communications between entities in a wide-area network. In a pub-sub system, *publishers* submit information to the system, while *subscribers* can register to

receive publications of interest. Data is routed through a network of *routing nodes* that form an overlay network in the system. Pub-sub systems scale to handle large volumes of data from many applications by establishing routing paths that efficiently deliver messages to subscribers while eliminating duplicate messages along those paths.

In this paper, we are particularly interested in pub-sub systems that support *in-network aggregation*, in which routing nodes perform hierarchical aggregation on data that is published to the system. In-network aggregation is useful for applications that monitor the state of wide-area control systems—such as the electric power grid [26] and building management systems (BMSs) [11]—by collecting individual readings from an array of sensors and other devices. The sensors and other devices monitoring the control system act as publishers that push measurements to the pub-sub system, while monitoring applications act as subscribers that later receive those events from the system. In some cases, sensors deployed over a wide physical area can generate data at very high rates. For example, phasor measurement units [22] used in the electric power grid generate data 30 times per second. Since most monitoring applications make control decisions based on aggregated data computed from raw sensor data, the demanding bandwidth and latency requirements of the pub-sub system can be reduced through the use of hierarchical in-network aggregation.

Safely supporting hierarchical in-network aggregation in a pub-sub system requires that we address two important security issues. First, publishers should be able to protect their individual inputs from potentially untrusted routers and subscribers. For example, in the power grid, utilities must hide their market sensitive input data from their competitors. Similarly, in BMS systems, the occupancy of a certain room in a building could reveal the occupants' private activities, while aggregate occupancy information for sections of the of the building is likely to be safe to disclose. Since routing nodes in a wide-area pub-sub system are typically managed by entities in different administrative domains, we must assume that publishers do not trust routing nodes in other domains in terms of the confidentiality of their published data. Therefore, our secure aggregation protocol should allow untrusted routing nodes to perform aggregation of raw data without learning the private input values.

Second, subscribers should have some assurance regarding the authenticity and integrity (correctness) of the aggregate data that they receive. Data integrity is critical for monitoring applications used in control systems, as these systems

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

WPES'08, October 27, 2008, Alexandria, Virginia, USA.

Copyright 2008 ACM 978-1-60558-289-4/08/10 ...\$5.00.

make control decisions based upon aggregated data received from the pub-sub system. Without this property, malicious routers could modify the data aggregated by the system, thereby tricking monitoring applications into making unsafe control decisions. Ensuring this integrity property on aggregated data is challenging because subscribers may not completely trust the routing infrastructure used by system. Furthermore, since thousands of publishers might contribute to a single aggregate data item, subscribers require a means of verifying the correctness of aggregate data without checking signed raw data, as this defeats the purpose of in-network aggregation.

In this paper, we present a secure aggregation protocol for the additive aggregate function *sum*. We believe that this is a reasonable first step towards the general secure aggregation problem for pub-sub systems because we can reduce other useful aggregation functions—such as *average*, *count*, *variance*, and *standard deviation*—from the *sum* function. Our protocol allows a collection of publishers to publish an aggregate result derived from their individual private inputs that is released only to authorized subscribers. Thus, such a pub-sub system enables users in the system to share useful statistical information without compromising the confidentiality of individual raw data. In addition, we apply a message authentication code (MAC) scheme based on the discrete logarithm property to allow subscribers to verify the correctness of aggregated results. Our scheme eliminates the necessity of providing a subscriber with the signed raw data used as input to the aggregation, and enables each subscriber to verify the correctness of aggregated data using an aggregated MAC of a constant size. We prove that our protocol satisfies the soundness and safety requirements of both publishers and subscribers.

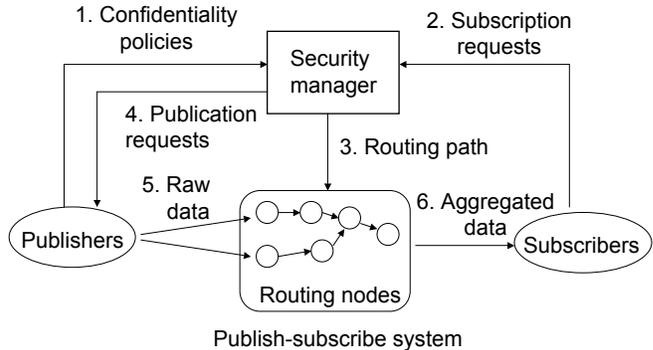
The rest of the paper is organized as follows. We describe our system and attack models for pub-sub systems in Section 2. We then present our aggregation protocol and proofs of its safety and security properties in Section 3. We discuss a possible way of extending our protocol to reduce our trust assumptions in Section 4. We cover related work in Section 5 and present our conclusions and directions for future work in Section 6.

## 2. SYSTEM MODEL

### 2.1 System overview

We assume that a pub-sub system consists of a set of routing nodes as well as a trusted security manager node, while publishers and subscribers exist as applications outside of the pub-sub system. We assume that every publisher, subscriber, and routing node is managed by some principal in the set  $\mathcal{P}$  of all principals. The security manager is a central authority that is trusted by publishers and subscribers to coordinate the system. Many existing pub-sub systems (e.g., [12, 19, 23]) that protect publishers’ private data requires such a central authority for key management. Each principal  $p_i \in \mathcal{P}$  maintains a public key pair  $(K_i, K_i^{-1})$ , and can obtain the public keys of other principals using a PKI or some other key distribution service. We assume that publishers publish data items from some value set  $\mathcal{V}$ .

Figure 1 describes the interactions between publishers, subscribers, and the components of the pub-sub system. First, each publisher notifies the security manager of the confidentiality policies protecting the aggregation of its data



**Figure 1: System model.** Each arrow is labeled with data transferred between two components.

values. As a result, the security manager has a complete view of all data confidentiality policies. When a principal wishes to begin a subscription, he issues a subscription request to the security manager. A subscription request is a set of (principal, variable) pairs in  $\mathcal{P} \times \mathcal{V}$  that represents the sum of variables that is to be computed. For example, a subscription request for the sum of publisher principal  $p_0$ ’s variable  $v_0$  and  $p_1$ ’s variable  $v_1$  is represented as  $\{(p_0, v_0), (p_1, v_1)\}$ . Publishers in our system model are more tightly coupled with subscribers than in many existing pub-sub systems; each subscriber explicitly specifies which publishers should provide raw data for the aggregate sum in a subscription request. However, this additional feature of a pub-sub system is essential to support monitoring applications for control systems because they need to use only data from trustworthy publishers to make right control decisions. Note that our pub-sub system can still choose an aggregation path while making it transparent from both publishers and subscribers.

After a subscription request is issued, the security manager checks whether the subscription request satisfies the confidentiality policies of the publishers owning the variables comprising the subscription request. If the appropriate confidentiality policies are satisfied, the security manager computes a hierarchical routing and aggregation path from the publishers to the subscriber, which is then sent to the appropriate routers. We assume that the security manager computes aggregation paths such that every routing node receives data from more than one publisher or routing node. As we will see in Section 3, our algorithms leverage this property of aggregation paths to ensure that each publisher’s data remains private.

Exactly how the security manager should compute this aggregation path is out of the scope of this paper—our secure aggregation protocol is independent of the aggregation path computation algorithm. At this point, the publishers can send their data to routing nodes, which forward the data while performing in-network aggregation. Finally, each subscriber receives the final aggregate data from the root routing node of the aggregation path.

Each publisher  $p_i$  publishes a variable  $v_i$  repeatedly synchronizing with the other publishers; that is, all the publishers publish their variables of each round simultaneously (i.e., with the same timestamp). We denote principal  $p_i$ ’s variable  $v_i$  at time  $t$  by  $v_i(t)$ . When a pub-sub system aggre-

gates multiple variables, it only aggregate variables with the same timestamp. Therefore, if each publisher  $p_i$  for  $i = 1$  to  $n$  publishes a variable  $v_i$ , a subscriber who subscribes to the sum of those variables will receive  $\sum_{i=1}^n v_i(t)$  at each time step  $t$ .

## 2.2 Confidentiality policies

Each publisher  $p_i$  can define confidentiality policies to limit access to the variables that it maintains. Specifically, a publisher  $p_i$  can define an access-control list (ACL)  $acl_i(v_i)$  to limit access to variable  $v_i$ . If  $acl_i(v_i) = \{p_j\}$ , only subscriber  $p_j$  can receive  $p_i$ 's variable  $v_i$ . If any other subscriber  $p_k$  issues a subscription request for the variable  $v_i$ , the security manager maintaining  $p_i$ 's confidentiality policies rejects  $p_k$ 's request.

Each publisher can also define an access-control list that protects the sum of multiple variables, some of which are maintained by other publishers. We define such an ACL by including multiple (principal, variable) pairs in the ACL. For example,  $p_i$  can define policy  $acl_i((p_i, v_i), (p_j, v_j))$  to protect the sum of  $p_i$ 's  $v_i$  and  $p_j$ 's  $v_j$ . When a subscriber  $p_k$  issues a subscription request  $\{(p_i, v_i), (p_j, v_j)\}$  for the sum of the two values,  $p_k$  must satisfy both  $p_i$  and  $p_j$ 's confidentiality policies on the sum of those two variables; that is,

$$\begin{aligned} p_k &\in acl_i((p_i, v_i), (p_j, v_j)) \text{ and} \\ p_k &\in acl_j((p_i, v_i), (p_j, v_j)) \end{aligned}$$

must hold. We can define confidentiality policies on the sum of more than two variables in the similar way.

## 2.3 Attack models

We now consider attacks on the system from the viewpoints of both publishers and subscribers. From the viewpoint of publishers, there are two kinds of adversaries. One is an adversary of colluding routing nodes that attempt to learn an aggregate sum in an unauthorized way. The other is an adversary who tries to learn some publisher's individual data value in an unauthorized way. The second type of adversary could include unauthorized subscribers as well as untrusted routing nodes. Such colluding parties can freely share messages that they obtain through the process of in-network aggregation. They can also determine an aggregation path by intercepting all the messages among all the parties who participate in the aggregation process. However, an adversary cannot see the contents of messages among non-colluding parties since all the messages between two parties are transmitted via secure channels established with pairwise shared keys. In this paper, when we consider an adversary who tries to learn individual data in an unauthorized way, we assume that the adversary consists of up to  $m - 1$  routing nodes and a subscriber; publishers cannot be part of the adversary because each publisher never receive information on the other publishers' confidential input data. We also assume that non-colluding parties follow our aggregation protocol properly and do not disclose information they obtain through the aggregation processes to the adversary.

In this paper, we do not address inference attacks carried out by subscribers with multiple subscriptions since this problem has been studied in the context of query auditing for statistical databases [6, 13, 14]. We assume that each publisher trusts the security manager to make the proper authorization decisions such that subscribers cannot infer unauthorized data from multiple subscriptions.

From the viewpoint of subscribers, attackers are untrusted routing nodes that incorrectly aggregate data values. Such an adversary can be comprised of any number of colluding routing nodes. On the other hand, we assume that publishers are not part of an adversary against subscribers because publishers can always modify the aggregated data by providing malicious input data while otherwise following the aggregation protocol properly. Therefore, if a subscriber  $p_{sub}$  subscribes to an aggregated result that includes some publisher  $p_i$ 's data,  $p_{sub}$  implicitly trusts  $p_i$  to provide correct input data. We note that each subscriber in our protocols in Section 3 must trust other subscribers who subscribe to the same subscription not to force it to accept incorrect aggregate values by colluding with malicious routing nodes. We will describe an extension of our protocol that removes this trust assumption in Section 4.

## 3. SECURE AGGREGATION IN A PUB-SUB SYSTEM

In this section, we describe our protocol incrementally. In Section 3.1, we first present a protocol that preserves publishers' confidentiality while aggregating data. We then present a protocol that ensures the integrity of aggregated data in Section 3.2. Finally, we present an integrated protocol that ensures both the confidentiality and integrity requirements of publishers and subscribers in Section 3.3.

### 3.1 Confidentiality-preserving aggregation

We now describe an aggregation protocol that preserves the confidentiality of each publisher's private input. Consider the case in which there are  $n$  publishers,  $p_1, \dots, p_n$ , and a single subscriber  $p_{sub}$  involved in the aggregation protocol. Each publisher  $p_i$  maintains a private variable  $v_i \in \mathbb{Z}$  and an associated confidentiality policy  $acl_i(v_i) = \emptyset$ . That is,  $p_i$  is not willing to disclose  $v_i$  to any other principal. However, every  $p_i$  is willing to disclose the sum  $\sum_{i=1}^n v_i$  to the subscriber  $p_{sub}$  and thus defines the confidentiality policy  $acl_i((p_1, v_1), \dots, (p_n, v_n)) = \{p_{sub}\}$ . The  $n$  publishers will use the pub-sub system to disseminate the sum of their private variables to the subscriber  $p_{sub}$ . Recall from Section 2.3 that we must consider attacks launched by coalitions of up to  $m$  colluding nodes, including router nodes and the subscriber  $p_{sub}$ .

Our protocol protects each principal  $p_i$ 's private data  $v_i$  from unauthorized routing nodes and subscriber  $p_{sub}$  by requiring that each  $p_i$  splits their value  $v_i$  into  $m$  shares, which are then sent to  $m$  different routing nodes. However, this scheme alone is insufficient for protecting the final aggregated sum from unauthorized routers, since eventually a single routing node will compute the final aggregation that will be distributed to  $p_{sub}$ . Therefore, we also require that each publisher  $p_i$  generates a random number  $q_i$  in  $\mathbb{Z}_p$  where  $p$  is a large integer and publishes  $m$  shares of the value  $v_i - q_i \pmod{p}$ . In our protocol, the subscriber  $p_{sub}$  will know the method that each  $p_i$  uses to choose  $q_i$  and thus can reconstruct the actual sum from the aggregate value that they eventually receive from the pub-sub system.

Our confidentiality-preserving aggregation protocol consists of the following steps:

**Initial secret sharing.** Each publisher  $p_i$  generates a secret  $q_i$  that is used in conjunction with a pseudo-random number generator. A PRNG generates a sequence

of unpredictable values given a seed; i.e.,  $PRNG : \mathbb{Z}_l \times \mathbb{N} \rightarrow \mathbb{Z}_p$ , where  $l$  is the key size and  $p$  is the output size. For convenience, we will define the PRNG as operating  $PRNG : \mathbb{Z}_l \times \mathcal{T} \rightarrow \mathbb{Z}_p$ , where  $\mathcal{T}$  is a set of all timestamps.

1. Each publisher generates a seed  $q_i$  randomly and sends  $q_i$  to the subscriber  $p_{sub}$  secretly.

Note that this secret sharing process need only be performed once at the beginning of a subscription request.

**Publication of data.** Publisher  $p_i$  publishes the value  $v_i$  at time  $t_i$  as follows:

1. Compute  $q_i(t_i) = PRNG(q_i, t_i)$
2. Compute  $v'_i = v_i - q_i(t_i)$ . This step is necessary to protect the sum  $\sum_{i=1}^n v_i$  from the untrusted routing node that computes the sum of all shares.
3. Randomly split  $v'_i$  into  $m$  shares  $v'_{i,1}, \dots, v'_{i,m}$  such that  $v'_i = \sum_{j=1}^m v'_{i,j}$ .
4. Send shares  $v'_{i,1}, \dots, v'_{i,m}$  to  $m$  different routing nodes.

**Aggregation on routing nodes.** Each routing node receives some number of shares from publishers or other routing nodes and sends the sum of these shares to the next routing node along the aggregation path. We assume that the security manager determines an aggregation path for each subscription request such that every routing node receives shares from more than one publisher or routing node. Each routing node performs the following steps:

1. Receive shares  $v_1, \dots, v_k$  from  $k$  publishers and/or routing nodes. Here  $k$  is the number of child nodes in the aggregation path. Note that each of these shares is associated with the same timestamp  $t_i$ .
2. Compute the sum of the shares  $v = \sum_{i=1}^k v_i$ .
3. Send  $(v, t_i)$  to the next routing node along the aggregation path specified by the security manager.

**Computation of the sum.** We assume that the security manager establishes an aggregation path such that there is a single routing node that computes the sum  $v'_{sum}$  of all the shares published by all of the publishers. After the subscriber  $p_{sub}$  receives the aggregate value  $v'_{sum}$  associated with timestamp  $t_i$ , it performs the following steps to compute the sum.

1. Compute  $q_i(t_i)$  as  $PRNG(q_i, t_i)$  for each  $i$ .
2. Compute sum  $v_{sum} = v'_{sum} + \sum_{i=1}^n q_i(t_i) = \sum_{i=1}^n v_i$ .

Note that a naive aggregation protocol has a communication overhead of  $O(n+r)$  where  $n$  is the number of publishers and  $r$  is the number of routing nodes. Our confidentiality-preserving aggregation protocol, however, has a communication overhead of the order  $O(nm+r)$  because each publisher must send  $m$  shares to the pub-sub system. We now make the following claims regarding the security properties of this protocol.

**THEOREM 1 (CONFIDENTIALITY OF AGGREGATE SUM).** *No coalition of colluding routing nodes can obtain the sum  $\sum_{i=1}^n v_i$ .*

**PROOF.** Note that the root routing node of an aggregation path can only obtain the sum  $v'_{sum} = \sum_{i=1}^n v'_i$  where  $v'_i = v_i - q_i(t_i)$ . Since  $v'_{sum}$  has no correlation with the actual sum  $v_{sum} = \sum_{i=1}^n v_i$ , it is impossible for the routing node to learn any information about  $v_{sum}$  without knowing a secret value  $q_i$  shared between each  $p_i$  and the subscriber  $p_{sub}$ .

Also, the routing node cannot learn any correlation among a sequence of aggregate sums, since the PRNG ensures that, even if previous  $q_i(t_i)$  values are known, it is infeasible to predict future values of  $q_i(t'_i)$ .  $\square$

**THEOREM 2 (CONFIDENTIALITY OF INDIVIDUAL DATA).** *Let  $m$  be the number of shares generated by each publisher. No colluding parties of up to size  $m$  that includes routing nodes and the subscriber  $p_{sub}$  can obtain any principal  $p_i$ 's private data  $v_i$ .*

**PROOF.** Without loss of generality, we discuss the confidentiality of the variable  $v_i$  maintained by publisher  $p_i$ . The same argument holds for the other publishers' variables. Publisher  $p_i$  splits  $v'_i = v_i - PRNG(q_i, t_i)$  into  $m$  shares and send each share to a different routing node. Since all of the messages sent between  $p_i$  and the routing nodes are encrypted using pairwise shared keys, the only way to restore  $v'_i$  is to obtain all the shares of  $v'_i$  from the  $m$  routing nodes receiving these shares directly from  $p_i$ .

Recall that the security manager constructs an aggregation path such that every routing node receives shares from more than one publisher or routing node. Therefore, the other routing nodes that do not receive those shares directly from  $p_i$  receive some value, which is a sum of multiple shares including those from principals other than  $p_i$ , and thus those nodes cannot extract a share or multiple shares of  $v'_i$  from the values they receive.

Restoring  $v_i$  from  $v'_i$  also requires the value  $q_i(t_i)$  shared between  $p_i$  and  $p_{sub}$ . Thus, at least  $m + 1$  colluding parties are required to obtain the value of  $v_i$ .  $\square$

## 3.2 Integrity-preserving aggregation

We next describe an aggregation protocol that ensures only the integrity of the computed aggregate sum. We describe the protocol using the same example aggregation scenario used in Section 3.1. Our integrity-preserving aggregation protocol allows the subscriber  $p_{sub}$  to verify the integrity of the sum by leveraging a homomorphic MAC scheme based on the discrete logarithm property. A publisher generates the MAC of a value  $v$  by computing  $MAC(v, g) = g^v$ , where  $g$  is a generator for some multiplicative group  $G_p$  of prime order  $p$ . We assume that every principal knows the large prime number  $p$  that parameterizes  $G_p$ , and that the generator  $g$  is a shared secret among the publishers and the subscriber  $p_{sub}$ . As we will see, the second assumption ensures that a malicious routing node cannot modify an aggregated sum without being detected by the subscriber  $p_{sub}$ . This MAC scheme has the homomorphic property  $MAC(v_1, g) \times MAC(v_2, g) = MAC(v_1 + v_2, g)$  since  $g^{v_1} g^{v_2} = g^{v_1 + v_2}$ . Our integrity-preserving aggregation protocol consists of the following steps:

**Initial secret sharing.** All of the publishers and the subscriber  $p_{sub}$  share a secret generator  $g$ , which is cho-

sen by the security manager. Therefore, the subscriber  $p_{sub}$  needs to trust the security manager not to disclose  $g$  to the routing nodes. We assume that publishers will not disclose  $g$  to the routing nodes since a malicious publisher can always modify the aggregated sum by providing malicious input data. Each publisher  $p_i$  and the subscriber  $p_{sub}$  also share a secret seed  $r_i$ , similar to the  $q_i$  used in the confidentiality-preserving protocol in Section 3.1.

1. The security manager generates a generator  $g$  for the group  $G_p$  randomly and sends  $g$  to each publisher  $p_i$  and subscriber  $p_{sub}$  secretly. It is assumed that all principals (including routers) know the prime order  $p$  of  $G_p$ .
2. Each publisher generates a seed  $r_i$  randomly and sends  $r_i$  to the subscriber  $p_{sub}$  secretly.

**Publication of data.** The publisher produces the published value as well as a MAC.

1. Compute  $r_i(t_i) = PRNG(r_i, t_i)$
2. Compute the MAC as  $c(v_i) = g^{v_i+r_i(t_i)}$
3. Send  $v_i, c(v_i)$  to a routing node

Note that the MAC value is “blinded” by the use of the secret value  $r_i(t_i)$ . This prevents a known-plaintext attack from recovering the value of  $g$ .

**Aggregation on routing nodes.** During the protocol, each routing node receives some number of (value, MAC) pairs from publishers and/or other routing nodes. The router then computes the sum of those values and the product of those MACs, and passes the results to the routing node along the aggregation path. Each routing node performs the following steps:

1. Receive the (value, MAC) pairs  $(v_1, c_1), \dots, (v_k, c_k)$  from other publishers and/or routing nodes in the system. Note that each pair is associated with the same timestamp  $t_l$ .
2. Compute the sum of the shares  $v = \sum_{i=1}^k v_i$ .
3. Compute the product of the MACs  $c = \prod_{i=1}^k c_i \pmod{p}$ .
4. Send  $(v, c)$  and the timestamp  $t_l$  to the next routing node along the aggregation path.

**Computation and verification of the sum.** After subscriber  $p_{sub}$  receives  $v_{sum}$  and the MAC  $c_{sum}$  from the last routing node in the pub-sub system,  $p_{sub}$  computes the sum  $v_{sum}$  and verifies its correctness using the MAC  $c_{sum}$  as follows:

1. Receive  $v_{sum}$  and  $c_{sum}$ .
2. Calculate  $r_i(t_i)$  as  $PRNG(r_i, t_i)$  for  $i = 1$  to  $n$ .
3. Accept sum  $v_{sum}$  if  $g^{v_{sum}+\sum_{i=1}^n r_i(t_i)} \pmod{p} = c_{sum}$  holds.

Our integrity-preserving protocol requires the same number of messages as a simple aggregation protocol that does not ensure the integrity of the aggregated data. Although each message of our protocol must also include a MAC of

1024 bits, each subscriber can check the integrity of aggregate data only with a single MAC. On the other hand, if we take a naive approach of sending all the raw data and its signatures to each subscriber, the subscriber can check the integrity of aggregate data by recomputing that aggregate data from the signed raw data. However, this naive approach requires each subscriber to receive a message whose size is proportional to the number of publishers, and each subscriber will thus become a bottleneck of the system in terms of both computation and communication overhead.

We now formally prove that our protocol ensures the integrity of aggregated data.

**THEOREM 3 (SOUNDNESS).** *The probability that the subscriber  $p_{sub}$  accepts an incorrect  $v_{sum} \neq \sum_{i=1}^n v_i$  is no more than  $\frac{1}{p}$  where  $p$  is the prime order of group  $G_p$ .*

**PROOF.** We consider an adversary who controls *all* of the routing nodes in the pub-sub system. The inputs to this adversary would be the pairs  $v_i, c_i = g^{v_i+r_i(t)}$ . Consider a different generator  $\hat{g}$  in  $G_p$  where  $\hat{g} \neq g$ . Then there exist  $\hat{r}_i(t)$  such that  $g^{v_i+r_i(t)} = \hat{g}^{v_i+\hat{r}_i(t)}$ . Furthermore, assuming that the PRNG is cryptographically strong, the adversary has no way to distinguish between cases when  $g$  and  $\hat{g}$  are used. It is easy to see that:

$$c_{sum} = g^{v_{sum}+\sum_{i=1}^n r_i(t)} = \hat{g}^{v_{sum}+\sum_{i=1}^n \hat{r}_i(t)}$$

However, for any  $v'_{sum} \neq v_{sum}$ , it must be that:

$$g^{v'_{sum}+\sum_{i=1}^n r_i(t)} \neq \hat{g}^{v'_{sum}+\sum_{i=1}^n \hat{r}_i(t)}$$

since otherwise:

$$\begin{aligned} g^{v'_{sum}+\sum_{i=1}^n r_i(t)} / c_{sum} &= g^{v'_{sum}-v_{sum}} \\ = g^{v'_{sum}+\sum_{i=1}^n \hat{r}_i(t)} / c_{sum} &= \hat{g}^{v'_{sum}-v_{sum}}, \end{aligned}$$

which contradicts the above assumption  $g \neq \hat{g}$ . Therefore, for any  $v'_{sum} \neq v_{sum}$ , the adversary can only calculate the correct MAC and have it accepted by guessing the choice of  $g$ . Since it cannot do this without compromising the PRNG, the adversary can at best guess randomly, resulting in a  $\frac{1}{p}$  chance of success.  $\square$

### 3.3 Secure aggregation

We now present a secure aggregation protocol that ensures the integrity of aggregated data while also preserving the confidentiality of both individual input values and the aggregated sum. We develop this new protocol by combining the two previous protocols presented in Sections 3.1 and 3.2. In this protocol, each publisher  $p_i$  publishes shares of the MAC of variable  $v_i$  as well as shares of the value obtained by subtracting some random number from  $v_i$ . The protocol consists of the following steps:

**Initial secret sharing.** In this step, the publishers share two seeds,  $q_i$  and  $r_i$ , with the subscriber, as described in the privacy-preserving and integrity-preserving protocols.

**Publication of data.** In this step, each publisher  $p_i$  publishes  $m$  shares that sum to the value  $v_i - PRNG(q_i, t_l)$  and  $m$  other shares whose product is  $MAC(v_i + PRNG(r_i, t_l), g)$ . In particular, each publisher performs the following steps:

1. Compute  $q_i(t_i) = PRNG(q_i, t_i)$  and  $r_i(t_i) = PRNG(r_i, t_i)$
2. Compute  $v'_i = v_i - q_i(t_i)$ .
3. Split  $v'_i$  into  $m$  shares  $v'_{i,1}, \dots, v'_{i,m}$  randomly such that  $v'_i = \sum_{j=1}^m v'_{i,j}$ .
4. Compute  $v''_i = v_i + r_i(t_i)$ .
5. Split  $v''_i$  into  $m$  shares  $v''_{i,1}, \dots, v''_{i,m}$  randomly such that  $v''_i = \sum_{j=1}^m v''_{i,j}$ .
6. Compute the MAC  $c(v''_{i,j}) = MAC(v''_{i,j}, g)$  for  $j = 1$  to  $m$ .
7. Send the (value, MAC) pairs  $(v'_{i,1}, c(v''_{i,1})), \dots, (v'_{i,m}, c(v''_{i,m}))$  to  $m$  different routing nodes.

**Aggregation on routing nodes.** Each routing node receives multiple shares of values and MACs from publishers and/or other routing nodes. The router then computes the sum of the value shares and the product of the MAC shares. Each routing node performs the following:

1. Receive the (value, MAC) pairs  $(v_1, c_1), \dots, (v_k, c_k)$ , from some set of publishers and/or other routing nodes. Note that each (value, MAC) pair is associated with the same timestamp  $t_l$ ,
2. Compute the sum of the shares  $v = \sum_{i=1}^k v_i$ .
3. Compute the product of the MACs  $c = \prod_{i=1}^k c_i \pmod{p}$ .
4. Send  $(v, c)$  and the timestamp  $t_l$  to the next routing node along the aggregation path specified by the security manager.

**Computation and verification of the sum.** After the subscriber  $p_{sub}$  receives the sum  $v'_{sum}$  and the MAC  $c_{sum}$  from the last routing node in the aggregation path, they perform the following steps to compute the sum of  $v_1, \dots, v_n$  and verify its correctness:

1. Compute  $r_i(t_i) = PRNG(r_i, t_i)$  and  $q_i(t_i) = PRNG(q_i, t_i)$  for each  $i$ .
2. Compute the sum  $v_{sum} = v'_{sum} + \sum_{i=1}^n q_i(t_i) = \sum_{i=1}^n v_i$ .
3. Accept the sum  $v_{sum}$  if  $g^{v_{sum} + \sum_{i=1}^n r_i(t_i)} = c_{sum}$  holds.

This secure aggregation protocol has the same communication complexity as the protocol presented in Section 3.1, and each message includes a MAC value of 1024 bits to ensure the integrity of aggregate data.

We now consider the security properties of this secure aggregation protocol. The secure aggregation protocol preserves the confidentiality of the aggregate sum and individual data items as in Theorems 1 and 2. The proofs for this protocol are the same as those for the confidentiality-preserving protocol in Section 3.1, and thus we omit them here. Likewise, soundness is ensured by an argument similar to Theorem 3.

## 4. SECURE AGGREGATION WITH AN AUTHENTICATED MAC

In the secure aggregation protocol described in Section 3.3, the generator  $g$  used in the MAC scheme must be kept secret from the untrusted routing nodes. Otherwise, a malicious router could modify a pair  $(v'_{sum}, c_{sum})$  to another valid pair  $(v'_{sum} + k, c_{sum} \times g^k)$  for any  $k \in \mathbb{N}$ , thereby forcing a subscriber to accept the incorrect aggregate value. Therefore, in that section, we made the assumption that each subscriber trusts the other subscribers not to collude with the untrusted routing nodes. In this section, we show how this assumption can be removed by taking a *delayed verification* approach, which is taken by TESLA [17], to ensuring the integrity of received aggregate values.

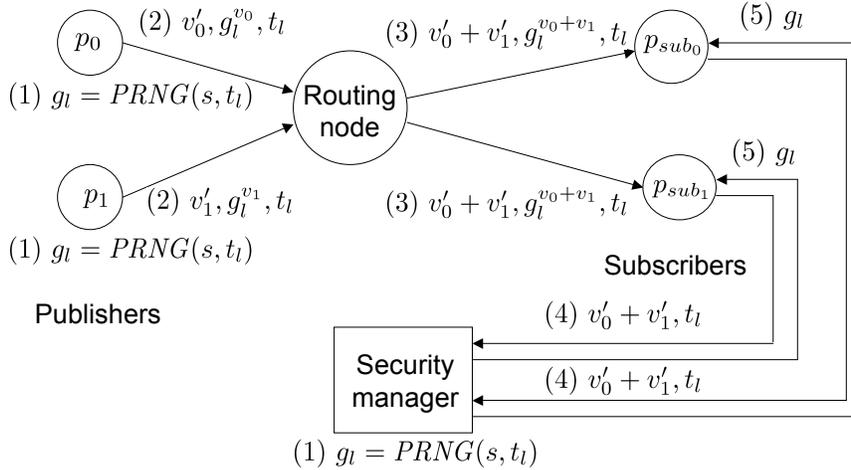
Figure 2 illustrates the mechanism used by our approach to providing subscribers with an authenticated MAC. For brevity, we only show a single routing node that sends the final result to the subscribers in Figure 2. In our scheme, the publishers and the security manager share a secret value  $s$  that is used to generate a pseudorandom sequence of generators for  $G_p$  that will be used by the MAC scheme. At each time  $t_l$ , a new generator  $g_l$  is created by the publishers and used to commit to values published during the current round of the aggregation. Value and MAC shares are then published and aggregated by the pub-sub system as described in Section 3.3, and eventually received by the subscribers in the system. After a subscriber receives a  $(v'_{sum}, c_{sum})$  pair at time  $t_l$ , he notifies the security manager that the value has been received. After the security manager receives such a notification from *every* publisher, it creates the generator  $g_l$  with the secret  $s$  and timestamp  $t_l$ . The generator  $g_l$  is then released to all of the publishers, who can then verify the integrity of the sum that was computed.

Notice that the subscribers cannot create the generator  $g_l$  themselves only from the timestamp  $t_l$ , as doing so requires knowledge of the secret  $s$ , which is shared only between the publishers and the security manager. As a result, *no* subscriber can verify the integrity of the aggregation that they receive during a particular round of the negotiation until *all* subscribers have received the aggregated value computed during that round. By forcing the system to proceed in a lockstep manner, the threat of malicious subscribers colluding with the routing nodes in the system to force other subscribers to accept incorrect aggregate values is clearly eliminated. This benefit comes at the cost of delaying the verification of the sums computed during the protocol; whether this delay is acceptable is application dependent, and may vary depending upon the environment in which our algorithm is deployed.

It is still possible for a malicious subscriber to perform a denial-of-service (DoS) attack by delaying sending a notification. However, the security manager can detect such subscribers and terminate their subscription requests if necessary.

## 5. RELATED WORK

In this section, we discuss related research efforts. We first examine existing security solutions for pub-sub systems and next discuss secure aggregation protocols for sensor networks. We finally cover research on verification of aggregated queries on outsourced databases briefly.



**Figure 2: Example of authenticating an aggregated MAC.** The arrows show the sequence of messages. For brevity, we only show the root routing node of a routing path in a pub-sub system.

### 5.1 Security in a pub-sub system

Many researchers have explored the security issues that arise in pub-sub systems. Wang et al. [28] describe various security issues with varying trust assumptions in a pub-sub system. In this paper, we focus on the issue of publication confidentiality and integrity under the presence of untrusted routing nodes, to which Wang et al. do not provide any concrete solutions.

Several researchers [29] have studied policy languages and enforcement mechanisms for limiting access to events in pub-sub systems with trusted routing nodes. Miklos [15] provides a policy language that defines access-control policies as filters in a pub-sub system. Zhao and Sturman [29] implement an access-control mechanism as a message filter on trusted routing nodes. Pesonen et al. [18] provides a scheme for delegation-based access control in a pub-sub system involving multiple security domains. Opyrchal et al. [16] develop an efficient key distribution scheme based on techniques of key caching to ensure confidentiality from end-point routing nodes to groups of subscribers. In this paper, we consider policy enforcement in systems with untrusted routing nodes.

Recently, researchers have begun to investigate security issues in pub-sub systems with untrusted routing nodes. However, none of this research addresses security issues associated with aggregated data. Khurana’s scheme [12] ensures publication confidentiality against routing nodes by encrypting confidential fields of each event with a key shared between a publisher and subscribers. Pesonen et al. [19] also use encryption on event attributes to protect confidential data from untrusted routing nodes. Raiciu et al. [20] developed several security protocols that allow routers to perform content-based filtering based on equality, keywords, and numeric values while keeping publications and subscriptions confidential from those routers. EventGuard [24] provides a comprehensive solution to various security problems in pub-sub systems that involve untrusted routing nodes. EventGuard ensures publication confidentiality by encrypting events with a shared key shared by publishers and subscribers. EventGuard has a central authority to issue a group key per each topic. EventGuard also ensures publication integrity by requiring publishers to sign their events.

However, EventGuard does not address publication confidentiality and integrity for aggregated data.

Ahmad et al. [1] developed a secure additive aggregation protocol in a large-scale overlay network. Their protocol uses an additively homomorphic public-key cryptosystem to protect confidential data from intermediate aggregation nodes. However, their scheme does not address the issue of integrity discussed in this paper. That is, there is no way for a subscriber to verify that no malicious intermediate node added an encrypted share multiple times to change the value of the sum.

### 5.2 Secure aggregation in a sensor network

Secure aggregation that ensures the confidentiality of raw data and the integrity of aggregated data has been mainly studied in the context of wireless sensor networks. In a sensor network, there is a single sink node that obtains aggregated data derived from sensor readings. Since the sink node is trusted to read all sensor readings, we do not consider the attack by colluding sensor nodes and the sink node, as we did for colluding parties of routing nodes and a subscriber.

Wagner [27] studies which aggregation functions can be securely computed in the presence of a few compromised nodes providing malicious input data. His model assumes that aggregation is performed in a single sink node of a sensor network. In this paper, we consider the case in which each subscriber trusts the publishers of aggregated data to provide correct inputs, and thus our security model excludes his threat model.

A number of aggregation protocols have been developed to ensure the confidentiality of sensor data from intermediate aggregator nodes. PDA [9] supports additive aggregation using a technique of secret splitting that is similar to ours, while protecting each sensor reading from other sensor nodes. However, PDA does not ensure the integrity of aggregated data. CDA [7] supports additive aggregation using an additively homomorphic encryption scheme to protect confidential sensor data from intermediate nodes that perform aggregation. In CDA, sensor nodes performing aggregation cannot report sensing data. Castelluccia et al. [4] also allows aggregation of encrypted data using an additively ho-

homomorphic stream cipher where each sensor node can have a different shared key with a sink node. Therefore, sensor nodes publishing data and aggregation nodes do not have to be disjoint. Solutions based on homomorphic encryption are not applicable to our problem since a malicious routing node that colludes with an unauthorized subscriber can forward an encrypted individual data to the subscriber who can decrypt it.

There are a few aggregation protocols that ensure the integrity of aggregated data. Hu and Evans [10] develop a secure aggregation protocol that ensures the integrity of aggregated data. In their protocol, each sensor node publishes data along with a message authentication code (MAC) constructed using a key generated every time it sends new data. Each node's parent forwards that data and its MAC along with the MAC for the aggregated data it receives to its parent. At every round of the protocol, a base station needs to broadcast the keys used by sensor nodes at the previous round so that each sensor node can verify the MACs that it received from other nodes. If a pair of child-parent nodes in the sensor network is compromised, their protocol cannot ensure the integrity of aggregated data. On the other hand, our scheme ensures integrity under the presence of an arbitrary number of untrusted routing nodes. Chan et al. [5]'s aggregation protocol ensures the integrity of additive aggregation by enabling sensor nodes to construct a commitment tree in a distributed fashion while they perform in-network aggregation. The base station receiving the sum verifies its correctness by sending the root vertex of the commitment tree to all sensor nodes using an authenticated broadcast. If each sensor node can verify that its contribution was added to the sum by checking the commitment tree, the base station accepts the sum. Although Chan's protocol does not involve expensive cryptographic operations, it is not applicable to a pub-sub system where routing nodes are deployed across a wide-area network because the authenticated broadcast would be inefficient in this case. The number of messages necessary for our integrity-preserving protocol in Section 3.2 is same as that of a simple aggregation protocol that does not preserve the integrity of aggregate data. On the other hand, Chan's protocol requires additional messages for an authenticated broadcast and for exchanging information (i.e., off-path values) on a commitment tree among sensor nodes. If the commitment tree contains  $n$  leaf sensor nodes, the sensor nodes need to exchange  $O(n \log n)$  messages in addition to the cost of forwarding those messages via intermediate sensor nodes.

### 5.3 Verification of aggregated queries

Haber et al. [8] address the problem of verifying the integrity of aggregate queries on outsourced databases. They develop a verification protocol that allows a user to verify the integrity of the sum of multiple values in a database without seeing those individual values. This query is processed by an untrusted service provider that is different from the trusted database owner. Since a single database owner provides all of the individual values for the sum, their protocol can ensure the integrity of the sum by providing the user with a Merkle hash tree of commitments to the individual values, so that the user can verify the authenticity of those commitments with a digital signature on the root node of the hash tree created by the database owner. This solution of constructing a single digital signature on mul-

iple commitments is not applicable to our problem, since each individual data is provided by a different publisher.

## 6. CONCLUSION

In this paper, we presented a secure aggregation protocol for computing the additive function *sum* in a pub-sub system. Our protocol allows an aggregated value to be computed from the raw inputs of some number of publishers in a privacy-preserving manner. Secret splitting is used to ensure that as long as no more than  $m$  parties collude, no principal's private data value will be leaked. Our protocol further guarantees that the computed aggregate is disclosed only to authorized subscribers and cannot be inferred by the untrusted routing infrastructure that comprises the pub-sub system. In addition, our scheme allows subscribers to verify the correctness of the aggregate value computed by the system by leveraging a homomorphic message authentication (MAC) scheme based on the discrete logarithm property. This MAC scheme allows the correctness of an aggregation to be verified by subscribers and imposes a small constant-size data transmission overhead, regardless of the number of publishers contributing to the aggregate value computed by the system.

In the future, we hope to develop support for other useful aggregate functions, such as *min* and *max*. We also plan to incorporate a fault tolerance mechanism that allows partial aggregates to be computed in the presence of failed publisher nodes. Supporting such a mechanism will require an access control policy language that specifies both disclosure constraints, as well as constraints on the number of failed nodes that will be tolerated by a given publisher.

## Acknowledgments

We would like to thank Patrick Tsang for his comments on a draft of this paper, the anonymous reviewers for their helpful suggestions, and the entire Trustworthy Cyber Infrastructure for Power Grid project team for discussions at an early stage of this research. This research has been supported in part by grant CNS 05-24695 (CT-CS: Trustworthy Cyber Infrastructure for the Power Grid (TCIP)) from the US National Science Foundation.

## 7. REFERENCES

- [1] Waseem Ahmad and Ashfaq Khokhar. Secure aggregation in large scale overlay networks. *Proceedings of the 49th Global Telecommunications Conference*, pages 1–5, November 2006.
- [2] David E. Bakken, Carl H. Hauser, Harald Gjermundrod, and Anjan Bose. Towards more flexible and robust data delivery for monitoring and control of the electric power grid. Technical Report TR-GS-009, Washington State University, May 2007.
- [3] Antonio Carzaniga, David S. Rosenblum, and Alexander L. Wolf. [Design and evaluation of a wide-area event notification service](#). *ACM Transactions on Computer Systems*, 19(3):332–383, August 2001.
- [4] Claude Castelluccia, Einar Mykletun, and Gene Tsudik. Efficient aggregation of encrypted data in wireless sensor networks. In *The Second Annual Conference on Mobile and Ubiquitous Systems: Networking and Services*, pages 109–117, July 2005.

- [5] Haowen Chan, Adrian Perrig, and Dawn Song. Secure hierarchical in-network aggregation in sensor networks. In *Proceedings of the 13th ACM conference on Computer and communications security*, pages 278–287, New York, NY, USA, 2006. ACM.
- [6] Francis Chin. Security problems on inference control for sum, max, and min queries. *J. ACM*, 33(3):451–464, 1986.
- [7] Joao Girao, Markus Schneider, and Dirk Westhoff. On concealed data aggregation in wireless sensor networks. In *Proceedings of IEEE International Conference on Communication*, May 2005.
- [8] Stuart Haber, William Horne, Tomas Sander, and Danfeng Yao. Privacy-preserving verification of aggregate queries on outsourced databases. Technical Report HPL-2006-128, HP Labs, December 2006.
- [9] Wenbo He, Lue Liu, Hoang Nguyen, Klara Nahrstedt, and Tarek Abdelzaher. Pda: Privacy-preserving data aggregation in wireless sensor networks. *26th IEEE International Conference on Computer Communications*, pages 2045–2053, May 2007.
- [10] Lingxuan Hu and David Evans. Secure aggregation for wireless networks. In *Proceedings of the 2003 Symposium on Applications and the Internet Workshops*, page 384, Washington, DC, USA, 2003. IEEE Computer Society.
- [11] Wolfgang Kastner, Georg Neugschwandtner, Stefan Soucek, and Michael H. Newmann. Communication systems for building automation and control. *Proceedings of the IEEE*, 93(6):1178–1203, June 2005.
- [12] Himanshu Khurana. Scalable security and accounting services for content-based publish/subscribe systems. In *Proceedings of the 2005 ACM symposium on Applied computing*, pages 801–807, New York, NY, USA, 2005. ACM Press.
- [13] Francesco M. Malvestuto and Mauro Mezzini. Auditing sum queries. In *Proceedings of the 9th International Conference on Database Theory*, pages 126–142, London, UK, 2002. Springer-Verlag.
- [14] Francesco M. Malvestuto, Mauro Mezzini, and Marina Moscarini. Auditing sum-queries to make a statistical database secure. *ACM Transactions on Information System Security*, 9(1):31–60, 2006.
- [15] Zoltan Miklos. Towards an access control mechanism for wide-area publish/subscribe systems. In *Proceedings of the 22nd International Conference on Distributed Computing Systems*, pages 516–524, Washington, DC, USA, 2002. IEEE Computer Society.
- [16] Lukasz Opyrchal and Atul Prakash. Secure distribution of events in content-based publish subscribe systems. In *Proceedings of the 10th conference on USENIX Security Symposium*, pages 21–21, Berkeley, CA, USA, 2001. USENIX Association.
- [17] Adrian Perrig, Ran Canetti, Doug Tygar, and Dawn Song. Efficient authentication and signing of multicast streams over lossy channels. In *Proceedings of the 2000 IEEE Symposium on Security and Privacy*, pages 56–73, Washington, DC, USA, May 2000. IEEE Computer Society.
- [18] Lauri I. W. Pesonen, David M. Eyers, and Jean Bacon. A capability-based access control architecture for multi-domain publish/subscribe systems. In *Proceedings of the International Symposium on Applications on Internet*, pages 222–228, Washington, DC, USA, 2006. IEEE Computer Society.
- [19] Lauri I. W. Pesonen, David M. Eyers, and Jean Bacon. Encryption-enforced access control in dynamic multi-domain publish/subscribe networks. In *Proceedings of the 2007 inaugural international conference on Distributed event-based systems*, pages 104–115, New York, NY, USA, 2007. ACM.
- [20] Costin Raiciu and David S. Rosenblum. Enabling confidentiality in content-based publish/subscribe infrastructures. *Securecomm and Workshops*, pages 1–11, 2006.
- [21] Venugopalan Ramasubramanian, Ryan Peterson, and Emin Gun Sirer. Corona: A high performance publish-subscribe system for the world wide web. In *Proceedings of the 3rd Symposium on Networked Systems Design and Implementation*, May 2006.
- [22] Jr. Robert O. Burnett, Marc M. Butts, and Patrick S. Sterlina. Power system applications for phasor measurement units. *Computer Applications in Power, IEEE*, 7(1):8–13, 1994.
- [23] Mudhakar Srivatsa and Ling Liu. Securing publish-subscribe overlay services with eventguard. In *Proceedings of the 12th ACM conference on Computer and communications security*, pages 289–298, New York, NY, USA, 2005. ACM Press.
- [24] Mudhakar Srivatsa and Ling Liu. Secure event dissemination in publish-subscribe networks. In *Proceedings of the 27th International Conference on Distributed Computing Systems*, page 22, Washington, DC, USA, 2007. IEEE Computer Society.
- [25] Robert Strom, Guruduth Banavar, Tushar Chandra, Marc Kaplan, Kevan Miller, Bodhi Mukherjee, Daniel Sturman, and Michael Ward. [Gryphon: An information flow based approach to message brokering](#). In *International Symposium on Software Reliability Engineering (ISSRE '98)*, November 1998.
- [26] Kevin Tomsovic, David E. Bakken, Vaithianathan Venkatasubramanian, and Anjan Bose. Designing the next generation of real-time control, communication, and computations for large power systems. *Proceedings OF THE IEEE*, 93(5):965–979, 2005.
- [27] David Wagner. Resilient aggregation in sensor networks. In *Proceedings of the 2nd ACM workshop on Security of ad hoc and sensor networks*, pages 78–87, New York, NY, USA, 2004. ACM.
- [28] Chenxi Wang, Antonio Carzaniga, David Evans, and Alexander L. Wolf. [Security issues and requirements for Internet-scale publish-subscribe systems](#). In *Proceedings of the 35th Annual Hawaii International Conference on System Sciences*, Big Island, Hawaii, January 2002.
- [29] Yuanyuan Zhao and Daniel C. Sturman. Dynamic access control in a content-based publish/subscribe system with delivery guarantees. In *Proceedings of the 26th IEEE International Conference on Distributed Computing Systems*, page 60, Washington, DC, USA, 2006. IEEE Computer Society.