

---

# Secure Data Aggregation Protocol for Sensor Networks

Presented By : Kavit Shah

Major Advisor : Brian King

Thesis Committee Members : Paul Salama, Mohamed El-Sharkawy, Sangkook Lee

---

# Overview

- Introduction
- Data Aggregation
- Secure Hierarchical In-network Data Aggregation (SHIA)
- Our Protocol
- Conclusion

---

# Introduction

**Ubiquitous Computing** - is a scenario in which computing is everywhere.

**Internet of Things** - is a system where the Internet is connected to physical world via ubiquitous sensors.

**Ad-hoc Networking** - is a local network of sensors formed by peer-to-peer communications.

**Sensor Networks** - Collectively, we refer these concept as Sensor Networks.

---

## Sensor Networks

In sensor networks, thousands of sensors may interact with each other and collect raw data.

The data is processed by a computationally powerful machine (the base station).

Then the base station converts data into information.

Based on the information an important action is taken.

---

## Sensor Networks Applications

**Military** - enemy tracking, battle field surveillance or target classification.

**Environmental** - to monitor geographical location without much human intervention.

**Health Care** - to monitor patients around the clock, send reminders to doctors and nurses.

**Sustainable Mobility** - to build digitally connected and coordinated vehicles.

---

## Consequences of Sensor Failures

Speed sensor failure led to crash of Air France flight - Airbus A330-203 AF 447 on 1st June 2009.

**For nearly a minute**, as the speed sensors jumped, the pilot was not present in the cockpit.

The pilots could not reclaim control as the plane dropped out of the sky at a rate of **10,000 feet per minute**.

The flight plunged into the Atlantic nose-up after its three-and-a-half minute freefall , killing all **228** on board.

---

## Resource Constrains in Sensor Network

**Physical Limitations** - often deployed in open, hostile and unattended environments.

**Hardware Limitations** - due to lower manufacturing cost of sensor nodes, they have low speed processor, limited storage, a short range trans receivers.

**Transmission Medium** - wireless networks have approximately  $10^6$  **times higher bit error rate (BER)** than wired networks which causes frequent link loss and then path loss.

**Mobility** - network topology is dynamic, topology changes due to link failure, node failure or bandwidth optimization.

---

# Cryptographic Tools

## Hash Functions

A hash function takes a message as its input and outputs a fixed length message called hash code, creating a digital fingerprint of the message.

**Compression** A function  $h$  maps an input  $x$  of arbitrary finite bitlength, to an output  $h(x)$  of fixed bitlength  $n$ .

**Preimage resistance** For all pre-specified outputs, it is computationally infeasible to find any input which hashes to that output, i.e., to find any preimage  $x'$  such that  $h(x') = y$  where  $y$  is given whose corresponding input is not known.

**2nd-preimage resistance** It is computationally infeasible to find any second input which has the same output as any specified input, i.e, given  $x$ , to find a 2nd-preimage  $x' \neq x$  such that  $h(x') = h(x)$ .

**Collision resistance** It is computationally to find any two distinct inputs  $x, x'$  which hash to the same output, i.e., such that  $h(x) = h(x')$ .

SHA-256, is a 256-bit hash and provides 128 bits of security against collision attacks.



---

## Message Authentication Codes

A Message Authentication Code (MAC) is a family of hash functions parameterized by a secret key  $k$ .

**Ease of computation** For a known function  $h_k$ , given a value  $k$  and an input  $x$ ,  $h_k(x)$  is easy to compute.

**Compression** A function  $h_k$  maps an input  $x$  of arbitrary finite bitlength to an output  $h_k(x)$  of fixed length  $n$ .

**Computation-resistance** Given a description of the function family  $h$ , for every fixed allowable value of  $k$  (unknown to an adversary), given zero or more text-MAC pairs  $(x_i, h_k(x_i))$ , it is computationally infeasible to compute any text-MAC pair  $(x, h_k(x))$  for any new input  $x \neq x'$  (including possibly for  $h_k(x) = h_k(x_i)$  for some  $i$ ). If computation-resistance does not hold, a MAC algorithm is subject to MAC-forgery.

---

## Digital Signatures

A digital signature is a cryptographic scheme for demonstrating the authenticity of a digital message.

A valid digital signature gives a recipient strong reason to believe that the message was created by a known sender, such that the sender can not deny having sent the message (**authentication and non-repudiation**) and that the message was not altered in transit (**integrity**).

It is a 5-tuple scheme  $(\mathcal{M}, \mathcal{K}, \mathcal{G}, \mathcal{S}, \mathcal{V})$

For any secret key  $s_k \in \mathcal{K}$  and any  $m \in \mathcal{M}$ , the message  $m$  is signed using key  $s_k$  as follows:

$$s = \text{Sign}_{s_k}(m) \quad (1)$$

For any  $s_k$  let  $p_k$  denote public key and for all  $m \in \mathcal{M}$  and  $s \in \mathcal{S}$ ,  $s$  as follows:

$$\text{Verify}_{p_k}(m, s) = \begin{cases} \text{true with probability of 1} & \text{if } s = \text{Sign}_{s_k}(m) \\ \text{false with overwhelming probability} & \text{if } s \neq \text{Sign}_{s_k}(m) \end{cases} \quad (2)$$

where the probability space is determined by the  $\mathcal{M}, \mathcal{S}, \mathcal{K}, \mathcal{K}'$  and perhaps the signing and verification algorithms. The “overwhelming probability” for the signature scheme determines the probability that the scheme allows for a forgery.

---

## Summary of Cryptography Tools

Three different integrity-protection mechanisms HASH, MAC, Signature can be summarized in a matrix like Table.

	Who can generate it	Who can verify it
Hash	Everyone	Everyone
MAC	Holders of secret	Holders of secret
Signature	Holder of secret	Everyone

---

## Data Aggregation

The idea of Data Aggregation is to compress the data coming from different sources enroute eliminating redundancy, minimizing the number of transmissions, thus saving energy and increasing the longevity of the network.

This paradigm shifts the focus from the traditional **address centric** approaches for networking (finding short routes between pairs of addressable end nodes) to a more **data-centric** approach (finding routes from multiple sources to a single destination that allows in-network consolidation of redundant data).

---

## Definition

Consider a sensor network with  $n$  sensors collecting data and the base station, processes data into information.

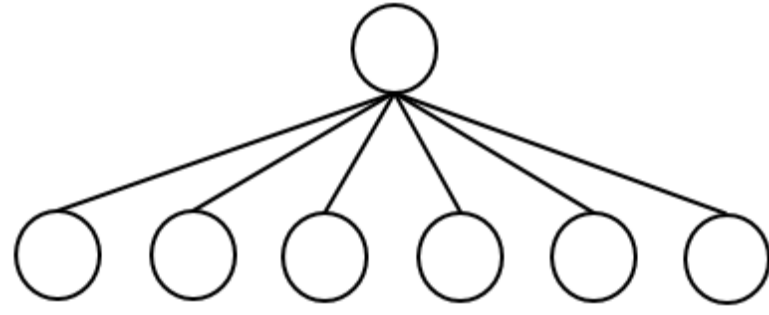
This can be represented as  $f(x_1, x_2, \dots, x_n)$  where  $x_1, x_2, \dots, x_n$  represent the sensor readings.

The goal is to compute the information as follows :

$$y = f(x_1, x_2, \dots, x_n)$$

---

## Data Aggregation using SUM



The root in the network, receives the following sensor data

$S_1(10)$ ,  $S_2(14)$ ,  $S_3(12)$ ,  $S_4(15)$ ,  $S_5(11)$ ,  $S_6(17)$  and has its own data  $S_0(15)$ .

The root node aggregates these seven data and creates an aggregated result as follows:

$$S = \sum_{i=0}^6 S_i \quad (3)$$

Now, root has to send only one data to its parent instead of seven.

---

## Lossy Data Compression

Lossy data compression schemes produces a compressed file from which only an **approximation** to the original information can be recovered. Much higher compression ratios are possible.

Our aggregation protocol is **lossy data compression** as the base station receives an aggregated sensor data and can not recover the original sensor data.

---

## Secure Hierarchical In-network Data Aggregation

SHIA is a light weight protocol for providing data-integrity to in-network data aggregation scheme.

It is designed for general networks with single or multiple adversaries.

Our work enhances SHIA, by making it communication efficient, adds new security services to the protocol, achieves similar security goals with non-resilient aggregation functions and efficient ways of analyzing the protocol.



---

## Network Assumptions

A multi hop network with a set  $S = \{S_1, \dots, S_n\}$  of  $n$  nodes where all sensor nodes are alive and reachable.

The network is organized in a **rooted tree topology**.

The **trusted base station** resides outside of the network and the network and the network topology.

The base station can directly communicate with every node in the network using **authenticated broadcast**.

All the wireless communications is **peer-to-peer** and we do not consider the **local wireless broadcast**.

Each sensor node  $I$  has a unique **ID** and shares a unique secret symmetric key  $sk_I$  with the base station.

The sensor nodes are capable of doing **symmetric and asymmetric** key encryption and decryption.

---

## Attacker Model

We consider a model with **polynomially bounded** adversary (polynomial in terms of the security parameter), which has a complete control over some of the sensor nodes in the network.

We focus on **stealthy attacks**, where the goal of an adversary is to make the base station accept false aggregation results, which are significantly different from the true results determined by the measured values, while not being detected by the base station.

We do not consider **denial-of-service (DoS)** and various **passive attacks**.

---

## SUM Aggregate Algorithm

Here, the aggregate function  $f$  is summation meaning that we want to compute  $a_1 + a_2 + \dots + a_n$ , where  $a_i$  is the sensor reading of the node  $i$ .

The algorithm has the following three main phases.

**Query Dissemination**, initiates the aggregation process.

**Aggregate Commit**, initiates the commitment tree generation process.

**Result Checking**, initiates the distributed, interactive verification process.

---

# Aggregation Tree

We require the prior constructed rooted aggregation tree.

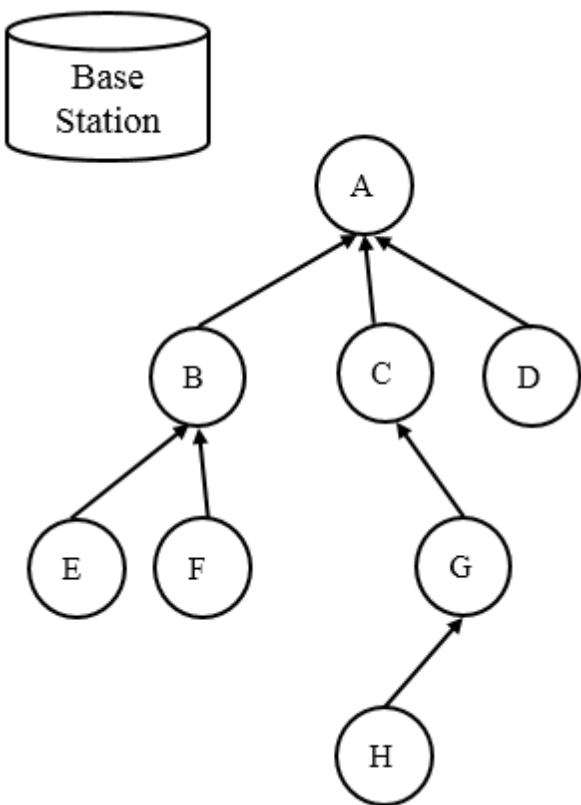


Figure 1: Aggregation Tree.

---

## Query Dissemination

The base station broadcasts the query request message with the query nonce  $N$  in the aggregation tree.

SHIA uses **hash chain** to generate new nonce for each query.

A hash chain is constructed by repeatedly evaluating a pre-image resistant hash function  $h$  on some initial random value, the final value (or “anchor value”) is preloaded on the nodes in the network.

The base station uses the pre-image of the last used value as the nonce for the next broadcast.



Figure 2: Hash Chain.

For example, if the last known value of the hash chain is  $h^i(X)$ , then the next broadcast uses  $h^{i-1}(X)$  as the nonce;  $X$  is the initial random value.

When a node receives a new nonce  $N'$ , it verifies that  $N'$  is a precursor to the most recently received (and authenticated) nonce  $N$  on the hash chain, i.e.,  $h^i(N') = N$  for some  $i$  bounded by a fixed  $k$  of number of hash applications.

---

## Data Item

A commitment tree is a binary tree where each vertex has an associated data-item representing the data that is passed on to its parent.

$$\langle id, count, value, commitment \rangle$$

Each sensor node creates its own data-item. For example, sensor node  $A$  creates its data-item  $A_0$ .

$$A_0 = \langle A_{id}, 1, A_{value}, H(N || 1 || A_{value}) \rangle$$

where  $A_{id}$ ,  $A_{value}$  is the unique ID and sensor reading of the node  $A$ . The count is 1 as there is only vertex in the subtree rooted at  $A$ ,  $H$  is the collision resistant hash function, and  $N$  is the query nonce.

---

## Signing and Verification of the Data-item

Each sensor node sends the signature of its data-item signed by itself using its own secret key.

$$S = \text{Sign}_{sk_A}(A_0)$$

Table: Digital Certificate

Unique ID of the sensor node
Public key of the sensor node
Certification Authority's name
Certification Authority's digital signature

$$\text{Verify}_{pk_A}(A_0, S) = \begin{cases} \text{true with probability of 1} & \text{if } S = \text{Sign}_{sk_A}(A_0) \\ \text{false with overwhelming probability} & \text{if } S \neq \text{Sign}_{sk_A}(A_0) \end{cases}$$

---

## Commitment Payload

A **commitment payload** is a set of data-items of the root vertices of the trees with their respective signatures in the outgoing commitment forest and an additional signature for the transmission.

The **transmit payload** is the concatenation of all the data-items in the commitment payload.



---

## Commitment Payload Example

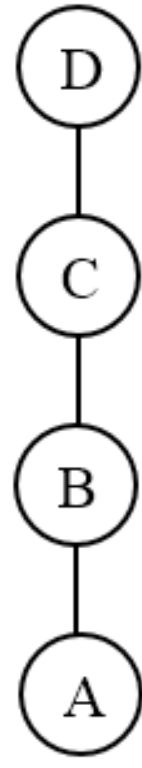


Figure 3: Palm Shaped Aggregation Tree

$$A_{pay} = \langle A_0, \text{Sign}_{sk_A}(A_0), \text{Sign}_{sk_A}(A_\tau) \rangle \text{ where } A_\tau = \langle A_0 \rangle$$

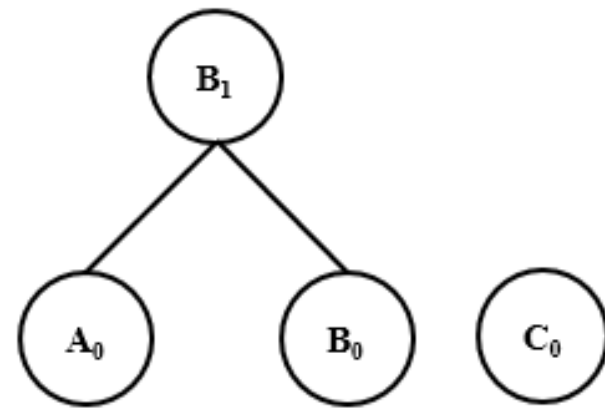


Figure 4: Commitment Payload of Sensor Node  $C$

$$C_{pay} = \langle C_0, \text{Sign}_{\text{sk}_C}(C_0), B_1, \text{Sign}_{\text{sk}_B}(B_1), \text{Sign}_{\text{sk}_C}(C_\tau) \rangle \text{ where } C_\tau = \langle C_0 || B_1 \rangle$$

$$C_0 = \langle C_{id}, 1, C_{value}, H(N || 1 || C_{value}) \rangle$$

$$B_1 = \langle B_{id}, 2, B_{1value}, H(N || 2 || B_{1value} || A_0 || B_0) \rangle; B_{1value} = B_{value} + A_{value}$$

---

## FSwRD vs FSwoRD

Forwarding Signatures With Resigning Data-Item (FSwRD)

$$C_{pay} = \langle C_0, \text{Sign}_{sk_C}(C_0), B_1, \text{Sign}_{sk_C}(B_1), \text{Sign}_{sk_C}(C_\tau) \rangle \text{ where } C_\tau = \langle C_0 || B_1 \rangle$$

Forwarding Signatures Without Resigning Data-Item (FSwoRD)

$$C_{pay} = \langle C_0, \text{Sign}_{sk_C}(C_0), B_1, \text{Sign}_{sk_B}(B_1), \text{Sign}_{sk_C}(C_\tau) \rangle \text{ where } C_\tau = \langle C_0 || B_1 \rangle$$

---

**Analogy for FSwRD vs FSwoRD**



Figure 5: Diamond Supply Chain.

---

## Security Benefits of Signatures

The signature allows the parent node to verify the **authenticity** of the sensor node and assures the **integrity** of the data-item.

It allows the sender to have the proof for the sent data-item and the receiver to have the proof for the received data-item, providing the security service of **non-repudiation**.

The digital signature depends on the message so the parent node can not reuse the signature for other messages in the future, protecting the network against the **replay attacks**.

The signature of the transmit-payload is like the signature for the transmission, assuring none of the data-items in its payload have been left stranded.

---

## Aggregate Commit

This phase creates the commitment tree for the given aggregation tree.

### Commitment Tree Generation

Leaf nodes in the aggregation tree construct and send their payload to their parents in the aggregation tree.

Each internal node in the aggregation tree constructs their leaf vertex.

Internal node verifies all the received signatures then it merges all the data-items with same count value from its forest.

It merges two data-items by creating a new data-item with count value incremented by one and whose value is the addition of value field of the previous two data-items.

For example, in Figure 1, the root  $A$  receives payloads from each of its children.

We describe the payload generation process for nodes  $D$ ,  $B$ ,  $C$  and  $A$  in order.

---

## Example Continue: D's Payload Generation

The node  $D$  constructs its payload and sends  $D_{pay}$  to its parent  $A$ .

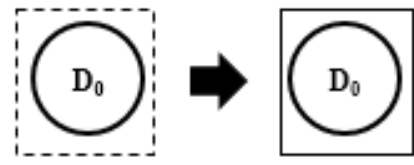


Figure 6: Transformation from  $D$ 's forest to its payload. Each dashed-line box shows forest and solid-line box shows payload of the respective sensor node.

$$D_{pay} = \langle D_0, \text{Sign}_{\text{sk}_D}(D_0), \text{Sign}_{\text{sk}_D}(D_\tau) \rangle; \text{ where } D_\tau = \langle D_0 \rangle$$

$$D_0 = \langle D_{id}, 1, D_{value}, H(N||1||D_{value}) \rangle$$

---

## Example Continue: B's Payload Generation

The node  $B$  constructs its payload from its forest which consists of payloads received from  $E$  and  $F$ .

Then node  $B$  sends  $B_{pay}$  to its parent  $A$ .

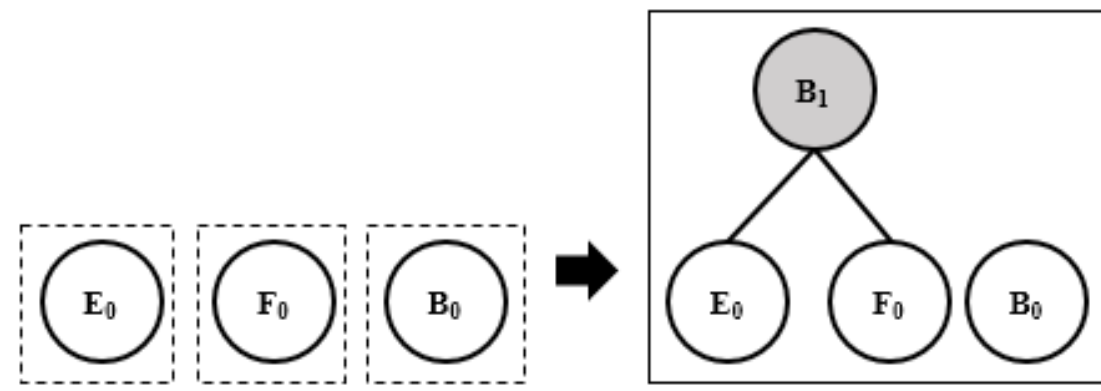


Figure 7: Transformation from  $B$ 's forest to its payload.

$$B_1 = \langle B_{id}, 2, B_{1value}, H(N||2||B_{1value}||E_0||F_0) \rangle; B_{1value} = E_{value} + F_{value}$$

$$B_{pay} = \langle B_0, \text{Sign}_{sk_B}(B_0), B_1, \text{Sign}_{sk_B}(B_1), \text{Sign}_{sk_B}(B_\tau) \rangle \text{ where } B_\tau = \langle B_0||B_1 \rangle$$



## Example Continue: C's Payload Generation

The node  $C$  constructs its payload from its forest which consists of payloads received from  $G$ .

Then node  $C$  sends  $C_{pay}$  to its parent  $A$ .

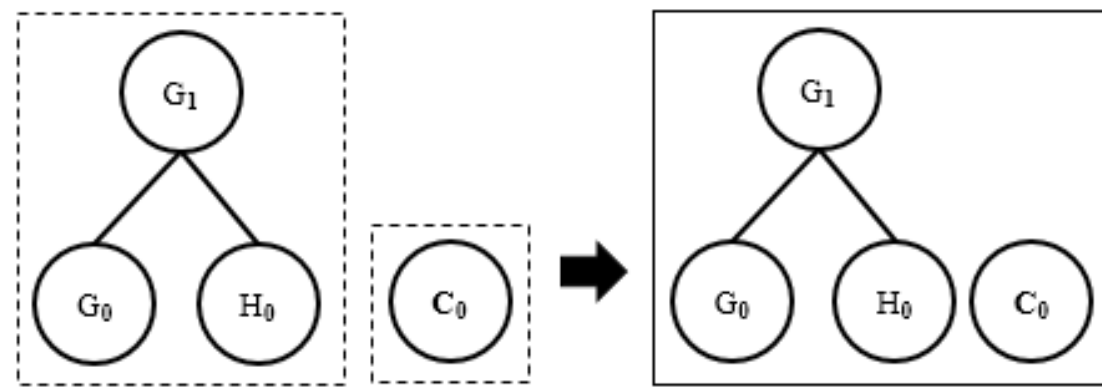


Figure 8:  $C$ 's forest aggregation creating its payload.

$$G_{pay} = \langle G_1, \text{Sign}_{sk_G}(G_1), \text{Sign}_{sk_G}(G_\tau) \rangle \text{ where } G_\tau = \langle G_0 || H_0 \rangle$$

$$C_{pay} = \langle C_0, \text{Sign}_{sk_C}(C_0), G_1, \text{Sign}_{sk_C}(G_1), \text{Sign}_{sk_C}(C_\tau) \rangle \text{ where } C_\tau = \langle C_0 || G_1 \rangle$$

---

**Example Continue: A’s Payload Generation**

The root node of the aggregation tree  $A$  receives the payloads from  $B$ ,  $C$  and  $D$  respectively.

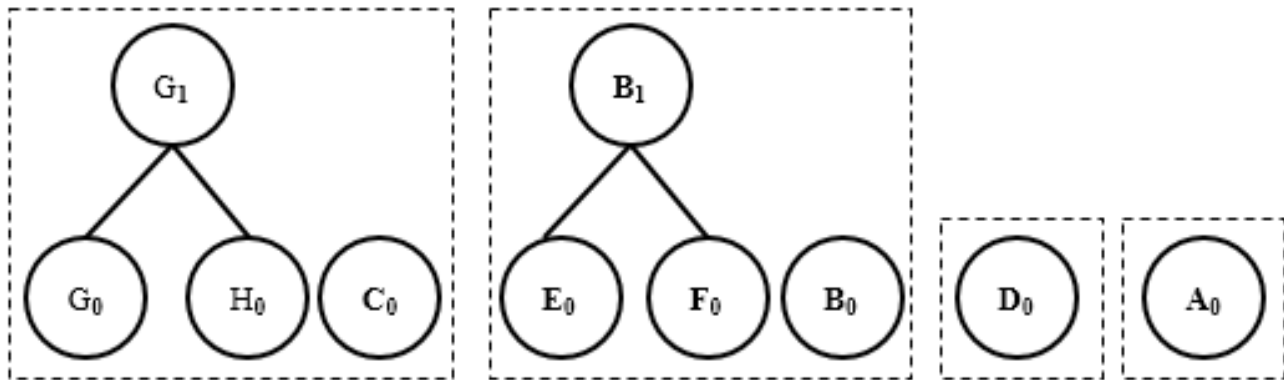


Figure 9:  $A$ ’s forest:  $A$  receives three payloads from  $C$ ,  $B$  and  $D$

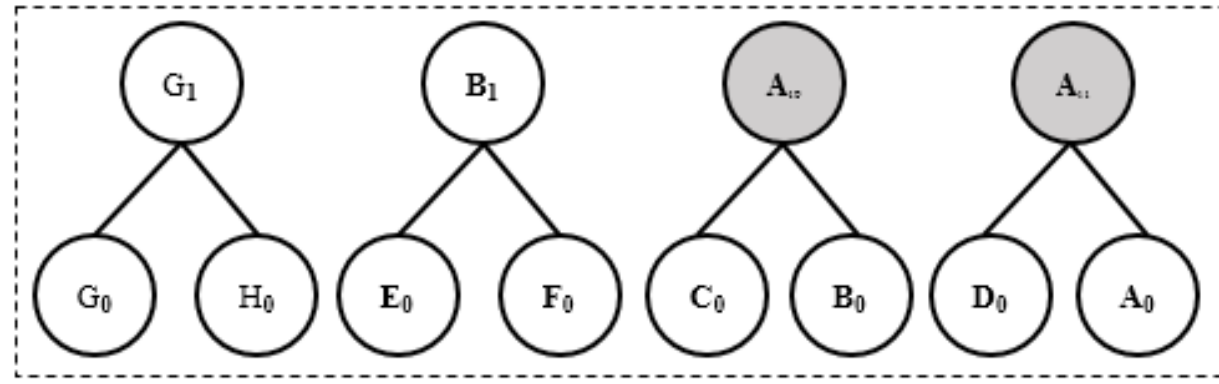


Figure 10:  $A$ 's forest: after first merge

$$A_{10} = \langle A_{id}, 2, A_{10value}, H(N||2||A_{10value}||B_0||C_0) \rangle; A_{10value} = B_{value} + C_{value}$$

$$A_{11} = \langle A_{id}, 2, A_{11value}, H(N||2||A_{11value}||D_0||A_0) \rangle; A_{11value} = D_{value} + A_{value}$$

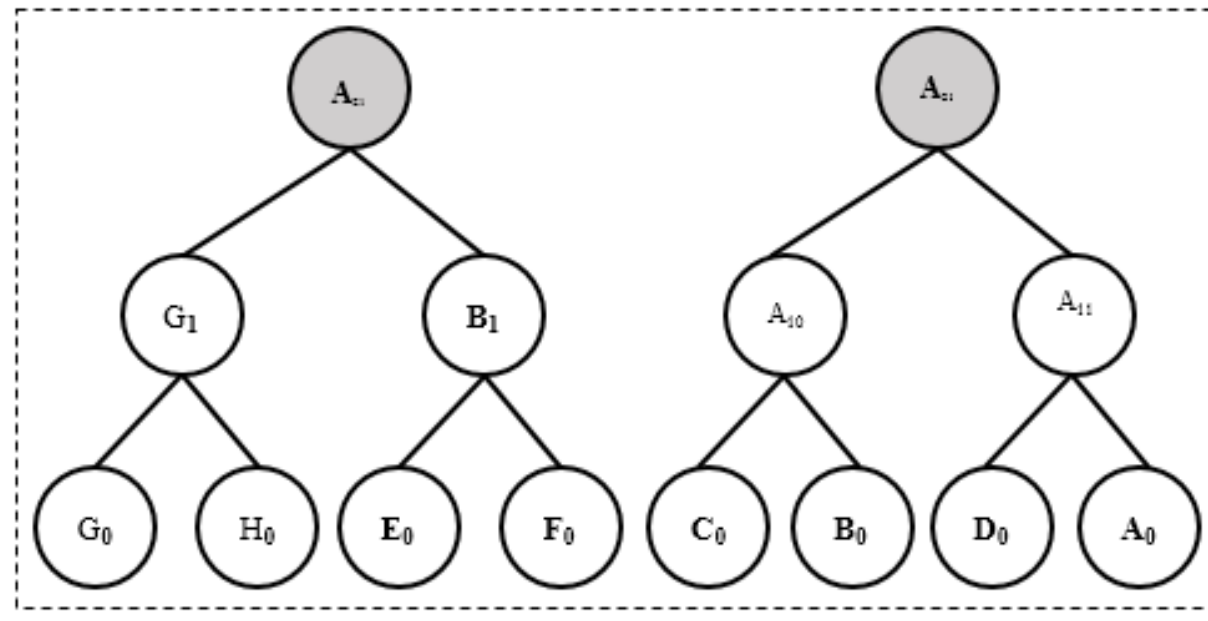


Figure 11:  $A$ 's forest: after second merge

$$A_{20} = \langle A_{id}, 4, A_{20value}, H(N||4||A_{20value}||G_1||B_1) \rangle; A_{20value} = G_{1value} + B_{1value}$$

$$A_{21} = \langle A_{id}, 4, A_{21value}, H(N||4||A_{21value}||A_{10}||A_{11}) \rangle; A_{21value} = A_{10value} + A_{11value}$$

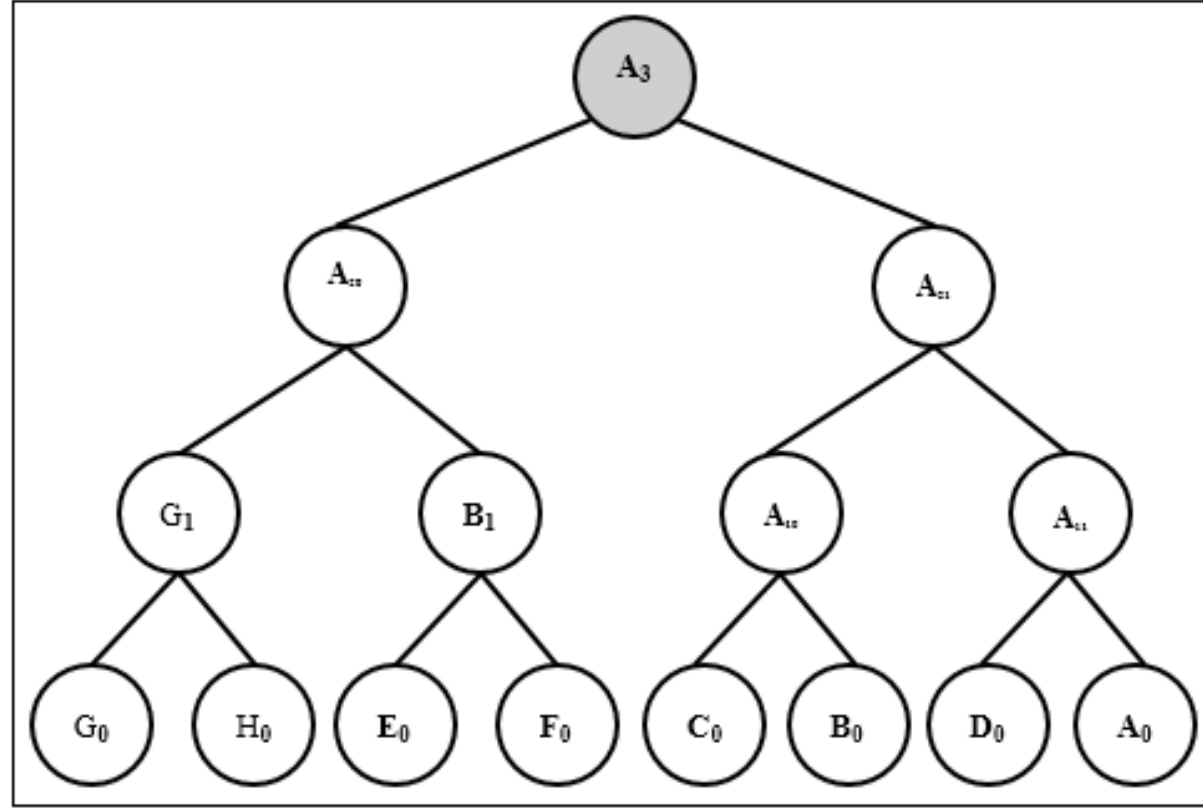


Figure 12:  $A$ 's payload :  $A$  sends  $A_3$  to the base station.

$$A_3 = \langle A_{id}, 8, A_{3value}, H(N || 8 || A_{3value} || A_{20} || A_{21}) \rangle; A_{3value} = A_{20value} + A_{21value}$$

$$A_{pay} = \langle A_3, \text{Sign}_{sk_A}(A_3), \text{Sign}_{sk_A}(A_\tau) \rangle \text{ where } A_\tau = \langle A_3 \rangle$$

---

## Result Checking

This phase requires that all the sensor nodes verify their individual contributions to the final aggregate value.

If there is any inconsistency in the aggregation process then with the help of the base station, trace down the node responsible for inconsistency.

It has the following major steps :

**Dissemination of Final Payload by the Base Station**

**Dissemination of Off-Path Values**

**Verification of Inclusion**

**Collection of Authentication Codes**

**Verification of Authentication Codes**

**Detecting An Adversary**

---

## Dissemination of Final Payload by the Base Station

The base station broadcasts all the data-items in the payload to entire network using **authenticated broadcast**.

The base station receives only one data-item  $A_3$  in the payload sent by  $A$ .

The base station broadcasts  $B_{pay}$  to entire network.

$$B_{pay} = \langle A_3, \text{Sign}_{sk_B}(A_3), \text{Sign}_{sk_B}(B_\tau) \rangle \text{ where } B_\tau = \langle A_3 \rangle .$$

---

## Dissemination of Off-Path Values

**Defn:** The set of **off-path vertices** for a vertex  $u$  in a tree is the set of all the siblings of each of the vertices on the path from  $u$  to the root of the tree that  $u$  is in (the path is inclusive of  $u$ ).

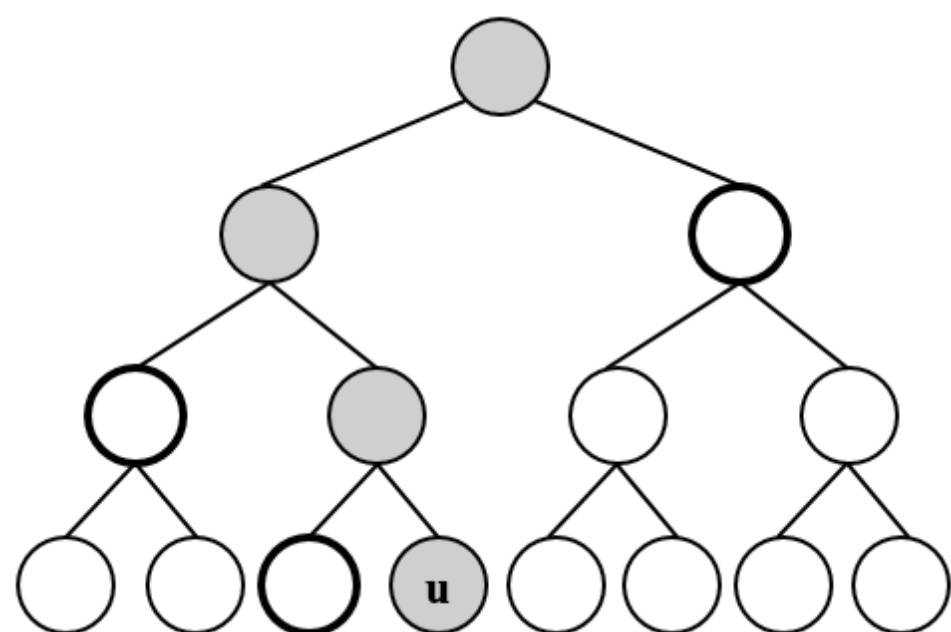


Figure 13: Off-path vertices of  $u$  are highlighted in bold.



---

In Figure 12,  $A_{10}$  has two children  $C_0$  and  $B_0$ .

$A_{10}$  also receives  $A_{11}$  and  $A_{20}$  from its parent  $A_{21}$ .

$A_{10}$  ( which is sensor node  $A$  in aggregation tree ) sends the following off-path values to  $C$  and  $B$  respectively.

$$< B_0, \text{Sign}_{sk_A}(B_0), A_{11}, \text{Sign}_{sk_A}(A_{11}), A_{20}, \text{Sign}_{sk_A}(A_{20}), \text{Sign}_{sk_A}(A_\tau) >$$

$$< C_0, \text{Sign}_{sk_A}(C_0), A_{11}, \text{Sign}_{sk_A}(A_{11}), A_{20}, \text{Sign}_{sk_A}(A_{20}), \text{Sign}_{sk_A}(A_\tau) >$$

---

## FSwoRD

$$< B_0, \text{Sign}_{sk_B}(B_0), A_{11}, \text{Sign}_{sk_A}(A_{11}), A_{20}, \text{Sign}_{sk_A}(A_{20}), \text{Sign}_{sk_A}(A_\tau) >$$

$$< C_0, \text{Sign}_{sk_C}(C_0), A_{11}, \text{Sign}_{sk_A}(A_{11}), A_{20}, \text{Sign}_{sk_A}(A_{20}), \text{Sign}_{sk_A}(A_\tau) >$$

In FSwoRD, all the leaf vertices need to know only one certificate as they receive data-items signed by their parent vertex.

In FSwoRD, all the leaf vertices might need to know  $\log l$  certificates, where  $l$  is the number of leaf-vertices in commitment tree.

---

## Significance of Commitment Filed

The commitment filed provides data-integrity and helps us detecting any *malicious activity* in the network. The signatures infrastructure eventually we can detect an adversary.

If an internal vertex simply **forwards incorrect data-item** then the relevant leaf vertex will complain, as they will not be able to derive the data-item received using authenticated broadcast from the base station.

If an internal vertex **changed the data-item** while creating commitment tree and sending the incorrect off-path values to compensate discrepancy.

---

## Malicious Activity

Suppose the vertices in the commitment tree have the data-items defined as follows :

$$A_0 = \langle A_{id}, 1, 10, H(N||1||10) \rangle$$

$$B_0 = \langle B_{id}, 1, 20, H(N||1||20) \rangle$$

$$C_1 = \langle C_{id}, 2, 30, H(N||2||30||A_0||B_0) \rangle$$

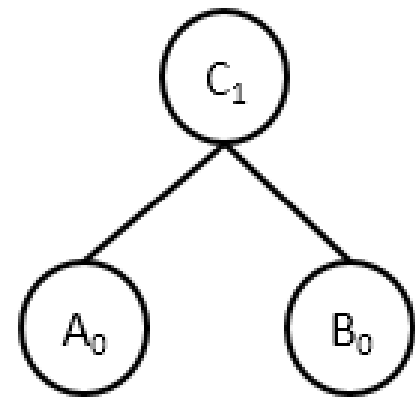


Figure 14: Smallest Possible Commitment Tree

---

Suppose  $C$  changes  $A_0$  and  $B_0$  to  $A'_0$  and  $B'_0$ .

$C$  can send either  $C'_1$  or  $C''_1$  to the base station.

To compensate for the discrepancy,  $C$  constructs  $B''_0$  and  $A''_0$  off-path values trying to hide its malicious activity from the base station.

$$A'_0 = \langle A_{id}, 1, 100, H(N||1||10) \rangle$$

$$B'_0 = \langle B_{id}, 1, 200, H(N||1||20) \rangle$$

$$C'_1 = \langle C_{id}, 2, 300, H(N||2||300||\mathbf{A''_0}||\mathbf{B_0}) \rangle \text{ or}$$

$$C''_1 = \langle C_{id}, 2, 300, H(N||2||300||\mathbf{A_0}||\mathbf{B''_0}) \rangle$$

$$B''_0 = \langle B_{id}, 1, 290, H(N||1||20) \rangle$$

$$A''_0 = \langle A_{id}, 1, 280, H(N||1||10) \rangle$$

---

$A$  and  $B$  receives either  $C'_1$  or  $C''_1$  from the base station based on what  $C$  has sent to base station.

$A$  and  $B$  receives  $B''_0$  and  $A''_0$  from  $C$  respectively.

$A$  and  $B$  derives the root data-item using the received off-path values, and it does not match with the received root data-item.

$A$  uses  $(A_0, B''_0)$  and derives  $\langle 2, 300, H(N||2||300||\mathbf{A}_0||\mathbf{B}''_0) \rangle = C''_1 \neq C'_1$

$B$  uses  $(A''_0, B_0)$  and derives  $\langle 2, 300, H(N||2||300||\mathbf{A}''_0||\mathbf{B}_0) \rangle = C'_1 \neq C''_1$ .

The commitment field makes it nearly impossible for an adversary to tamper with the data-items while creating commitment tree and/or while distributing off-path values.

