

Statistical Framework for Source Anonymity in Sensor Networks

Basel Alomair*, Andrew Clark*, Jorge Cuellar†, and Radha Poovendran*

*Network Security Lab (NSL), Electrical Engineering Department,
University of Washington, Seattle, Washington

†Siemens Corporate Technology, München, Germany
Email: {alomair,awclark,rp3}@uw.edu, jorge.cuellar@siemens.com

Abstract

In this work, we investigate the security of anonymous wireless sensor networks. To lay down the foundations of a formal framework, we propose a new model for analyzing and evaluating anonymity in sensor networks. The novelty of the proposed model is twofold: first, it introduces the notion of “interval indistinguishability” that is stronger than existing notions; second, it provides a quantitative measure to evaluate anonymity in sensor networks. The significance of the proposed model is that it captures a source of information leakage that cannot be captured using existing models. By analyzing current anonymous designs under the proposed model, we expose the source of information leakage that is undetectable by existing models and quantify the anonymity of current designs. Finally, we show how the proposed model can lead to a general and intuitive direction for improving the anonymity of current designs.

Index Terms

Wireless Sensor Networks (WSN), source location, privacy, anonymity

Submission category: regular paper

Contact author: Basel Alomair

1. Introduction

In sensor networks, small devices (called sensor nodes) are employed to capture relevant events and report collected data. The type of events nodes are designed to capture and report is an application dependent. Applications where sensor nodes can be utilized range from taking

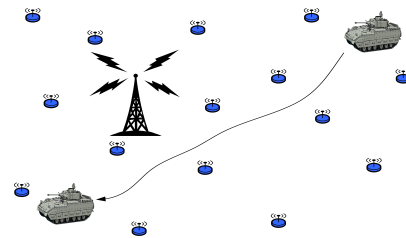


Figure 1. An example of a sensor networks deployed in a battlefield.

patients’ vital signs in controlled indoor environments to collecting tactical military information in hostile war zones.

A topic that has been drawing increasing research attention in wireless sensor networks is source location privacy [1]–[9]. (Source anonymity and source location privacy will be used synonymously for the rest of the paper.) Given the adversary’s knowledge of the locations of sensor nodes in the network, determining the individual nodes reporting the occurrence of real events can translate to the exposure of the location of the real events themselves. Applications in which hiding the occurrence of real events can be critical include, but are not limited to, the deployment of sensor nodes in battlefields as a means of coordinating strategic military actions, and the classic Panda-Hunter Game, where a malicious hunter monitors an existing animal tracking network to determine the location of the endangered panda [1], [2], [7], [8].

In such applications, at which source location privacy is of critical importance, special attention must be paid to the design of the node transmission algorithm so that monitoring sensor nodes does not reveal critical source information. One of the major challenges for the source anonymity problem is that it cannot be solved using tra-

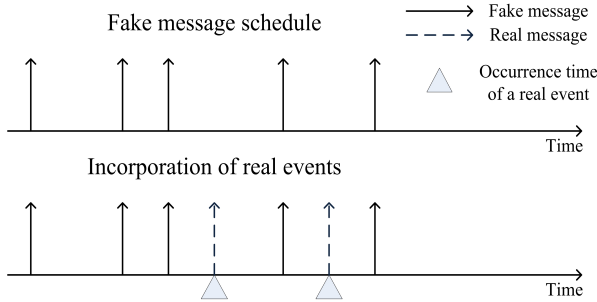


Figure 2. An illustration of the intuitive approach. The node is programmed to transmit fake messages so that real events are hidden within the fake transmissions.

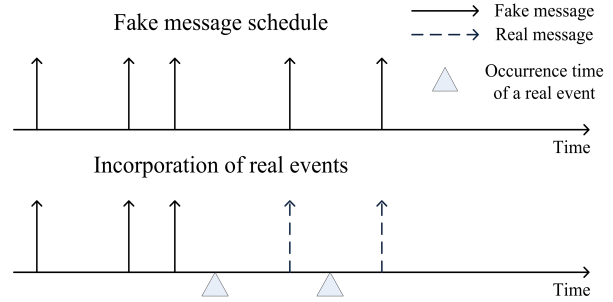


Figure 3. An illustration of the trivial solution. Nodes are programmed to transmit fake messages according to some pre-defined probabilistic distribution. When a real event occurs, it must be delayed until the next scheduled fake transmission before it can be sent out.

ditional cryptographic primitives. Encrypting nodes' transmissions, for instance, can hide the contents of plaintext messages, but the mere existence of ciphertexts is indicative of information transmission.

In the presence of a global adversary, who is able to monitor the traffic of the entire network, routing-based solutions has been shown to leak private source information [6]. An intuitive approach to report a real event without revealing, to a *global adversary*, its location information is to program nodes to transmit *fake messages* even if there are no real events to be reported [6]. When real events occur, they can be embedded within the transmissions of fake messages. This intuitive approach, however, does not completely solve the location privacy problem. When fake transmissions are scheduled according to some probabilistic distribution, statistical analysis can be used to distinguish between real and fake transmissions if real events are transmitted as they arrive. This intuitive approach is illustrated in Figure 2.

By realizing the problem with the intuitive approach of Figure 2, the solution becomes trivial. As opposed to transmitting real events as they occur, they can be transmitted instead of the next scheduled fake one. For example, sensor nodes can be programmed to transmit an encrypted message every minute. If there is no event to report, the node transmits a fake message. If a real event occurs within a minute from the last transmission, it must be delayed until exactly one minute after the last transmission has passed. This algorithm, trivially, provides source anonymity since an adversary monitoring a node will observe one transmission every minute and, assuming the semantic security of the underlying encryption, the adversary has no means of distinguishing between fake and real events. Figure 3 depicts an example of this trivial solution.

The trivial solution, however, has a major drawback: reporting real events must be delayed until the next scheduled transmission. (In the above example, in which a

transmission is scheduled every minute, the average latency of transmitting real events will be half a minute.) When real events have time-sensitive information, this latency might be unacceptable.

Reducing the latency of transmitting real events by adopting a more frequent scheduling algorithm is impractical for most sensor network applications. This is mainly because sensor nodes are battery powered and, in many applications, are unchargeable (for example, they maybe deployed in an unreachable or hostile environment). Consequently, a more frequent scheduling algorithm can exhaust nodes' batteries rather quickly, rendering sensor nodes useless.

Furthermore, a transmission scheduling based on any pre-specified probabilistic distribution, not necessarily deterministic as in the above example, will suffer the same problem discussed above: slower rates lead to longer latencies and faster rates lead to shorter battery lives. Consequently, practical solutions are designed to achieve the objective of source anonymity under two main constraints: minimizing latency and maximizing the lifetime of sensors' batteries.

To make things even more complex, the arrival rate and distribution of real events can be time varying and unknown in advance. Therefore, in the trivial solution, no pre-specified probabilistic distribution for fake transmissions can satisfy both constraints for arbitrary time-variant distribution of real event arrivals.

The current state of the art in designing anonymous sensor networks works as follows. In the absence of real events, nodes are programmed to transmit independent identically distributed (iid) fake messages according to a certain distribution with a certain rate. However, unlike the trivial solution, real events are transmitted as soon as possible (earlier than the next pre-scheduled fake transmissions) under the following condition: the distribution

of the entire message transmissions (fake and real) of each node is “statistically” similar to the transmission of only fake messages. (Statistical similarity is achieved via the use of statistical goodness of fit tests¹ to determine the transmission time of real events.) That is, to a global adversary monitoring the network, the time between any two transmissions (real or fake) will follow the same distribution of fake messages only. The current consensus is that this approach provides dependable solutions for the source anonymity problem in wireless sensor networks [7]–[11].

The above solution is better illustrated through the following example. Consider designing nodes to transmit fake messages according to a specific distribution. Further, let nodes store a sliding window of times between consecutive transmissions (inter-transmission times), say $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_k$, where \mathbf{X}_i is the random variable representing the inter-transmission time between the i^{th} and the $i+1^{\text{st}}$ transmissions and k is the length of the sliding window. Now, when a real event occurs, its inter-transmission time, denoted by \mathbf{X}_{k+1} , is defined to be the smallest value such that the sequence $\mathbf{X}_2, \mathbf{X}_3, \dots, \mathbf{X}_{k+1}$ passes some statistical goodness of fit tests. That is, an adversary observing the sequence of inter-transmission times will observe a sequence that is statistically indistinguishable from an iid sequence of random variables following the pre-specified distribution of fake transmissions.

However, continuing in this fashion will skew the mean since nodes always favor shorter inter-transmission times to transmit real events. To adjust the mean, the next inter-transmission time following a real event, \mathbf{X}_{k+2} in this example, must be stretched out. Again, \mathbf{X}_{k+2} is determined so that the sequence of inter-transmission times in the sliding window, $\mathbf{X}_3, \mathbf{X}_4, \dots, \mathbf{X}_{k+2}$, satisfies the same statistical goodness of fit tests used to compute \mathbf{X}_{k+1} . Therefore, an adversary observing the sensor node cannot distinguish between real and fake transmissions [7]–[11]. Figure 4 illustrates an instance of this approach.

In this paper, we take a closer look at the current state of the art in designing anonymous sensor networks. The driving motive behind this work is the key observation that, although an adversary might not be able to distinguish between real and fake transmissions, there still exists a source of information leakage that can affect the security of such designs. The inability to detect the source of information leakage in the current approach is not a result of false statements claimed in previous proposals; the lack of a formal framework that properly models anonymity in wireless sensor networks is the main reason for the inability to detect such a vulnerability. The main purpose

1. A statistical goodness of fit test is a statistical test that determines if a sequence of data samples follows a certain probabilistic distribution.

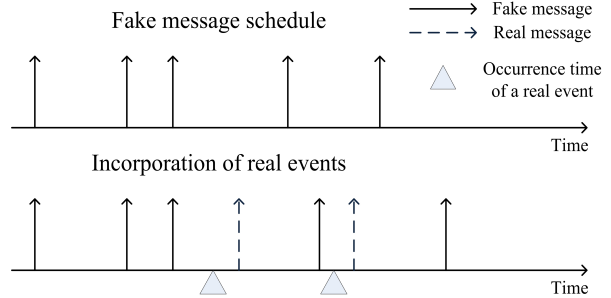


Figure 4. An illustration of the current state of the art in designing anonymous systems. Nodes transmit fake messages according to a pre-specified probabilistic distribution and maintain a sliding window of inter-transmission times. When a real event occurs, it is transmitted as soon as possible under the condition that the samples in the sliding window maintain the designed distribution. The transmission following the real transmission is delayed to maintain the mean of the distribution of inter-transmission times in the sliding window.

of this work is to provide such a framework.

1.1. Our Contributions

We summarize our contributions by the following points.

- We detect a source of information leakage in the current designs of anonymous systems that can undermine their security.
- We introduce the new notion of interval indistinguishability to analyze anonymity in wireless sensor networks. The new notion is stronger than existing notions and captures the source of information leakage that is undetectable by existing notions.
- We propose a quantitative measure to evaluate anonymity in sensor networks.
- We analyze, both analytically and via simulation, the current approach of designing anonymous sensor networks and quantify the amount of information leakage when the current approach is analyzed under the proposed model.
- Inspired by the new model, we suggest an approach for improving the anonymity of existing designs.

We emphasize that the goal of this work is not to propose a specific design for anonymous sensor network. This work aims to provide a general, security oriented, model for analyzing and evaluating the security of anonymous systems.

The rest of the paper is organized as follows. In Section 2 we discuss some related works. In Section 3 we provide

a formal definition for the proposed model and discuss how it properly captures the anonymity problem. Section 4 is devoted to the analysis of the current approach using the proposed model. In Section 5, we suggest a new direction for solving the anonymity problem in wireless sensor networks. We conclude the paper in Section 6.

2. Related Work

Source location privacy in sensor networks is part of a broader area, the design of anonymous communication systems. The foundation for this field was laid by Chaum in [12], and since then has become a very active area of research. In particular, topics related to location anonymity have been discussed by Reed *et al.* in [13], who introduced the idea of preserving anonymity through onion routing, and by Gruteser and Grunwald in [14], who discussed ways to provide anonymity in location-based services, such as Global Positioning Systems.

In wireless sensor networks, much of the work in source location privacy assumes a passive, local eavesdropper operating close to the base station. Privacy is maintained in such models through anonymous routing. The location privacy problem was first introduced in [1], [2]. The local eavesdropper model was introduced and the authors demonstrated that existing routing methods were insufficient to provide location privacy in this environment. They also proposed a phantom flooding scheme to solve the problem. In [3], Xi *et al.* proposed a new random walk routing method that reduces energy consumption at the cost of increased delivery time. Path confusion has also been proposed as an anonymity-preserving routing scheme by Hoh and Gruteser in [4]. In [5], Ouyang *et al.* developed a scheme in which cycles are introduced at various points in the route, potentially trapping the adversary in a loop and forcing the adversary to waste extra resources.

However, in the global adversarial model, in which the adversary has access to all transmissions in the network, routing-based schemes are insufficient to provide location privacy [6]. The global adversarial model was first introduced by Mehta *et al.* in [6]. The authors motivated the problem, analyzed the security of existing routing-based schemes under the new model, and proposed two new schemes. In the first scheme, some sensor nodes act as fake sources by mimicking the behavior of real events. For example, if the network is deployed to track an animal, the fake sources could send fake messages with a distribution resembling that of the animal's movements. This assumes some knowledge of the time distribution of real events, an assumption we do not make. In the second scheme, packets (real and fake) are sent either at constant intervals or according to a predetermined probabilistic schedule. Although this scheme provides perfect location privacy, it

also introduces undesirable performance characteristics, in the form of either relatively high latency or relatively high communication and computational overhead. The scheme of [7] was proposed to address this latency/overhead trade-off.

In [7], Shao *et al.* introduced the notion of *statistically strong source anonymity* in which a global adversary, who is able to monitor the traffic in the entire network, is unable to infer source locations by performing statistical analysis on the observed traffic. In order to realize their notion of statistical anonymity, nodes are programmed to transmit fake events according to pre-specified distribution. More specifically, after the transmission of every fake event, the node draws an exponentially distributed random variable $t \sim \text{Exp}(\lambda)$, where λ is the pre-specified rate of the exponential distribution. The node then waits for t time units and then transmits another fake event. That is, in the absence of real event transmissions, an adversary monitoring the sensor node will observe inter-transmission times that are iid exponentials with mean $\mu = 1/\lambda$.

Upon the occurrence of real events, the goal of a sensor node is to transmit them while maintaining the exponential distribution of the inter-transmission times. Obviously, if nodes delay their transmission of real events to the next scheduled fake transmission, no statistical test can be used to distinguish between real and fake events (since inter-transmission times are kept exponential iid's with the same rate). The goal, however, is to minimize the latency of reporting real events while maintaining statistical indistinguishability between real and fake transmissions.

To reduce the latency, the authors of [7] proposed the following procedure: let imd_i represent the inter-transmission time between the i^{th} and the $i + 1^{st}$ transmissions. Assume a real event has occurred after the transmission of the i^{th} event. Given $\{imd_1, imd_2, \dots, imd_i\}$, imd_{i+1} , the time after the transmission of the i^{th} event the node must wait before it can transmit the real event, is determined as follows: imd_{i+1} is the smallest positive value such that the sequence $\{imd_1, imd_2, \dots, imd_i, imd_{i+1}\}$ passes the Anderson-Darling (A-D) goodness of fit test [15] for a sequence of iid exponentials with mean μ .

Observe, however, that on average $imd_{i+1} < \mu$ since imd_{i+1} is, by definition, the minimum value that passes the test. Therefore, continuing in this fashion will cause the mean of the entire sequence to skew away from desired mean.

To solve the problem of mean deviation described above, the scheme in [7] includes a mean recovery algorithm. The mean recovery algorithm outputs a delay δ and the time between the transmission of a real event and the following event (fake or real) is set to $imd_{i+2} = t + \delta$, where $t \sim \text{Exponential}(\lambda)$. The scheme in [7] is designed so that the sequence $\{imd_1, \dots, imd_n\}$, where n is the

last transmitted message, always passes the A-D goodness of fit test. The authors also published a complementary paper [8] that minimizes the overall communication overhead by having some nodes act as proxies that filter out fake messages. This approach makes schemes based on generating fake messages more attractive by mitigating the performance problems.

Shao *et al.* also consider the problem of an active adversary in [16]. Their adversary also has the ability to perform node compromise attacks, and they develop tools to prevent the adversary from gaining access to event data stored in a node even if the adversary possesses that node's secret keys. Li *et al.* presented a survey of anonymity in sensor networks [9].

3. Modeling Anonymity

In this section we introduce our anonymity model for wireless sensor networks. Intuitively, anonymity should be measured by the amount of information about sources' locations an adversary can infer by monitoring the sensor network. The challenge, however, is to come up with an appropriate model that captures all possible sources of information leakage and a proper way of quantifying anonymity in different systems.

We start here by stating our assumptions about the network structure and the adversary's capabilities. We will then describe the currently used notion for source anonymity in sensor networks and point out a source of information leakage that is undetectable by this notion. Then, we will give a formal definition of a stronger anonymity notion that, in addition to capturing the sources of information leakage captured by the current notion, captures the source of information leakage that was missed by the current notion. Finally, we propose a quantitative measure for evaluating the security of anonymous sensor networks.

3.1. Network Model

We assume that communications take place in a network of energy constrained sensor nodes. That is, nodes are assumed to be powered with unchargeable batteries, thus, conserving nodes energy is a design requirement. Nodes are also equipped with a semantically secure encryption algorithm,² so that computationally bounded adversaries are unable to distinguish between real and fake transmission by means of cryptographic tests. When a node detects an

2. In cryptography, semantic security implies that, given a ciphertext, unauthorized users without the knowledge of the decryption key have no means of distinguishing between two plaintexts in which one of them corresponds to the observed ciphertext [17].

event, it places information about the event in a message and broadcasts an encrypted version of the message.

3.2. Adversarial Model

Our adversary is similar to the one considered in [6], [7], in that it is *external*, *passive*, and *global*. By external, we mean that the adversary does not control any of the nodes in the network and also has no control over the real event process. By passive, we mean that the adversary is capable of eavesdropping on the network, active attacks are not considered. By global, we mean that the adversary can simultaneously monitor the activity of all nodes in the network. In particular, the adversary can observe the timing and origin of every transmitted message.

As opposed to a global adversary, a local adversary is only capable of eavesdropping over a small area, typically the area surrounding the base station, and attempts to determine the source of traffic by examining the packet routing information or trying to follow the packets back to their source. Protocols that attempt to disguise the source of traffic through routing, while highly secure against local adversaries, do not defend against global adversaries [6].

We also assume that the adversary is capable of storing a large amount of message traffic data and performing complex statistical tests. Furthermore, the adversary is assumed to know the distribution of fake message transmissions. The only information unknown to the adversary is the timing when real events occur.

3.3. Event Indistinguishability (EI)

Currently, anonymity in sensor networks is modeled by the adversary's ability to distinguish between individual real and fake transmissions by means of statistical tests. That is, given a series of nodes' transmissions, the adversary should not be able to distinguish, with significant confidence, which transmission carries real information and which transmissions is fake.

Consider an adversary observing the sensor network over multiple time intervals, without being able to distinguish between individual fake and real nodes' transmissions. Assume, however, that during a certain time interval the adversary is able to notice a change in the statistical behavior of transmission times of a certain node in the network. This distinguishable change in transmission behavior can be indicative of the existence of real activities reported by that node, even though the adversary was unable to distinguish between individual transmissions.

For example, consider a sensor network deployed in a battlefield. For a certain time interval, there was no activity in the vicinity of a sensor node the enemy is monitoring. Therefore, by design, the node has been transmitting fake

messages for the duration of that time interval. Assume now that a moving platoon is in the vicinity of this node and the node started to report location information about the moving platoon. The enemy does not need to distinguish between individual transmissions to infer the existence of the moving platoon. All that is needed is the ability to distinguish between the time interval when no real activity is reported and the time interval when the platoon is in the vicinity of the sensor node.

Consequently, in many real life applications, modeling source anonymity in sensor networks by the adversary's ability to distinguish between *individual event transmissions* is insufficient to guarantee location privacy. This fact calls for a stronger model to properly address source anonymity in sensor networks.

Before we proceed to the new anonymity model, we formally define the currently adopted notion to model anonymity in sensor networks, namely, event indistinguishability.

Definition 1 (Event Indistinguishability - 'EI'): Events reported by sensor nodes are said to be indistinguishable if the inter-transmission times between them cannot be distinguished with significant confidence by means of statistical tests.

3.4. Interval Indistinguishability (II)

The main goal of source location privacy systems is to hide the existence of real events. This implies that, an adversary observing a sensor node during different time intervals, at which some of the intervals include the transmission of real events and the others do not, should not be able to determine with significant confidence which of the intervals contain real traffic.

This leads to the notion of interval indistinguishability that will be essential for our anonymity formalization.

Definition 2 (Interval Indistinguishability - 'II'): Let I_F denotes a time interval with only fake event transmissions (call it the "fake interval"), and I_R denotes a time interval with some real event transmissions (call it the "real interval"). The two time intervals are said to be statistically indistinguishable if the distributions of inter-transmission times during these two intervals cannot be distinguished by means of statistical tests.

To model interval indistinguishability, we propose the following game between a challenger \mathcal{C} (the system designer) and an adversary \mathcal{A} .

Game 1 (Modeling Interval Indistinguishability):

- 1) \mathcal{C} draws a bit $b \in \{0, 1\}$ uniformly at random.
- 2) \mathcal{C} chooses two intervals I_0 and I_1 , in which I_b is a real interval and the other one is fake.
- 3) \mathcal{C} gives I_0 and I_1 to \mathcal{A} .

- 4) \mathcal{A} makes any statistical test of her choice on I_0 and I_1 and outputs a bit b' .
- 5) If $b' = b$, \mathcal{A} wins the game.

With Definition 2 and Game 1, we aim to find a security measure that can formally quantify the anonymity of different designs. Let $\Pr[b' = b]$ be the adversary's probability of winning Game 1 and identifying the real interval. We quantify the anonymity of the sensor network by

$$\Lambda := 1 - 2\left(\Pr[b' = b] - \frac{1}{2}\right). \quad (1)$$

Observe that, in the best case, the adversary cannot do better than a random guess, i.e., $\Pr[b' = b] = 1/2$ leading to $\Lambda = 1$ (absolute anonymity). In the worst case, $\Pr[b' = b] = 1$ leading to $\Lambda = 0$ (no anonymity). Any other probability of winning the game will give an anonymity measure in the interval $[0, 1]$. Therefore, the anonymity measure of equation (1) is well-defined.

Given Definitions 1 and 2, the relation between event indistinguishability (EI) and interval indistinguishability (II) is stated as follows.

Lemma 1: Interval Indistinguishability \Rightarrow Event Indistinguishability.

Proof: Assume there exist a system satisfying interval indistinguishability but does not satisfy event indistinguishability. Since real and fake transmissions are distinguishable, given a fake interval and a real interval, the real interval can be identified as the one with the real transmission; a contradiction to the hypothesis that the system satisfies interval indistinguishability. Therefore, if intervals are indistinguishable, then individual events within them must also be indistinguishable. \square

To show that the proposed notion is stronger than the current one, it remains to show that event indistinguishability does not imply interval indistinguishability. Section 4 proves this fact by providing a counter example.

With the above definition of interval indistinguishability, we introduce the notion of Λ -anonymity in sensor networks.

Definition 3 (Λ -anonymity): A wireless sensor network is said to be Λ -anonymous if it satisfies two conditions

- 1) in different stages of each interval, inter-transmission times are indistinguishable,
- 2) the anonymity of the system, as defined in equation (1), is at least Λ .

The first condition in Definition 3 is different than event indistinguishability. It merely means that an adversary cannot identify the beginning, the middle, nor the end of any interval. It is necessary to ensure that there is no distinguishable transition region between intervals. If such a transition exists, it can lead to anonymity breach.

Table 1. A list of used terms and notations.

E_i	The random variable representing the event reported in the i^{th} transmission
X_i	The random variable representing the inter-transmission time between the i^{th} and the $i + 1^{st}$ transmissions
I_F	A fake interval: an interval consisting of fake events only
I_R	A real interval: an interval containing some real event transmissions
short inter-transmission times	inter-transmission times that are shorter than the mean of the pre-defined distribution
long inter-transmission times	inter-transmission times that are longer than the mean of the pre-defined distribution
short-long pattern	a short inter-transmission time followed by a long inter-transmission time

4. Analysis of EI-based Approaches

In this section we analyze, using our proposed model, systems that were shown to be secure under event indistinguishability; i.e., EI-based systems. We provide theoretical analysis showing that real and fake intervals can be statistically distinguishable. Then, we simulate an existing scheme to show that the simulation results coincide with the analytical results, and to quantify the anonymity of the simulated design. We start by a recapitulation of EI-based approaches for providing source anonymity in sensor networks.

4.1. EI-based Approaches

Recall that nodes are designed to transmit fake messages according to a pre-specified distribution. Furthermore, nodes store a sliding window of times between consecutive transmissions, say X_1, X_2, \dots, X_k , where X_i is the random variable representing the time between the i^{th} and the $i + 1^{st}$ transmissions and k is the length of the sliding window. When a real event occurs, its transmission time, represented by X_{k+1} , is defined to be the smallest value such that the sequence X_2, X_3, \dots, X_{k+1} passes some statistical goodness of fit tests. That is, an adversary observing the sequence of inter-transmission times will observe a sequence that is statistically indistinguishable from an iid sequence of random variables with the pre-specified distribution of fake message transmissions.

However, by continuing in this fashion, the mean will skew since nodes always favor shorter intervals to transmit real events. To adjust the mean, the next transmission following a real one, X_{k+2} in this example, will be delayed. Again, the delay is determined so that the sequence in the sliding window satisfies some statistical goodness of fit test. Consequently, as shown in [7], an adversary observing the sensor node cannot differentiate between real and fake transmissions.

4.2. Theoretical Interval Distinguishability

As discussed in Section 3, when an adversary can distinguish between real and fake intervals, source location

can be exposed, even if the adversary cannot distinguish between individual transmissions. In what follows, we give theoretical analysis of interval indistinguishability in EI-based systems.

Let X_i be the random variable representing the time between the i^{th} and the $i + 1^{st}$ transmissions and let $E[X_i] = \mu$. We will demonstrate an adversary's strategy of detecting the source location by investigating two intervals, a fake interval and a real one.

4.2.1. Fake Interval (I_F). In fake intervals, inter-transmission times are iid random variables. That is, the X_i 's are iid's with mean μ . Therefore, during any fake interval I_F , for any $X_{i-1}, X_i \in I_F$,

$$E[X_i | X_{i-1} < \mu] = \mu. \quad (2)$$

4.2.2. Real Interval (I_R). Let E_i be the random variable representing the event reported in the i^{th} transmission. Then, E_i can take the values R and F , where R denotes a real event transmission and F denotes a fake one. Since in general scenarios the distribution of inter-arrival times of real events can be varying and unknown beforehand, we will assume that E_i can take the values R and F with arbitrary probabilities.

Recall that the time between the transmission of a real event and its preceding fake one is usually shorter than the mean μ by design (to reduce latency). Recall further that the time between the transmission of a real event and its successive one is usually longer than μ by design (to adjust the ensemble mean). That is, during any real interval I_R , for any $X_{i-1}, X_i \in I_R$,

$$E[X_i | X_{i-1} < \mu, E_i = R] > \mu, \quad (3)$$

and,

$$E[X_i | X_{i-1} < \mu, E_i = F] = \mu, \quad (4)$$

by design. Using equations (3) and (4) we get,

$$\begin{aligned} & E[X_i | X_{i-1} < \mu] \\ &= E[X_i | X_{i-1} < \mu, E_i = R] \cdot \Pr[E_i = R] \\ &\quad + E[X_i | X_{i-1} < \mu, E_i = F] \cdot \Pr[E_i = F] \quad (5) \\ &> \mu \cdot \Pr[E_i = R] + \mu \cdot \Pr[E_i = F] \quad (6) \\ &= \mu. \quad (7) \end{aligned}$$

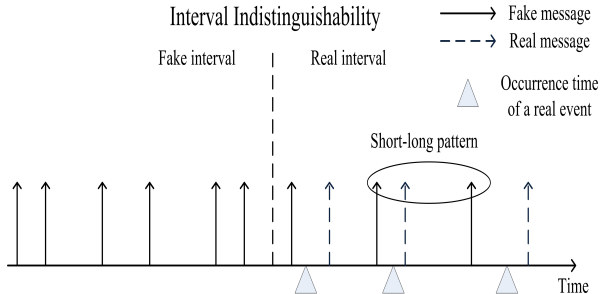


Figure 5. An illustration of interval distinguishability in the current approach. Real events are transmitted sooner than what is determined by the probabilistic distribution, while the transmission following the real event is later than what is determined by the probabilistic distribution to fix the mean of the pre-defined distribution.

Therefore, by equations (2) and (7), shorter inter-transmission times followed by longer inter-transmission times are most likely to occur in real intervals than fake intervals. This suggests the following strategy to distinguish between fake and real intervals: given two time intervals I_0 and I_1 , in which one of them is real and the other one is fake, the adversary counts the number of short followed by long inter-transmission times, simply called *short-long patterns* for the remainder of the paper. (An inter-transmission time is said to be short if its length is shorter than the mean μ , and is said to be long if it is longer than μ .) The interval that has more counts of short-long patterns is the real interval. Figure 5 illustrates the pattern of short-long inter-transmission times.

4.3. Case Study

In this section, we study the scheme appeared in [7], an instance of the EI-based approaches, and evaluate its anonymity using the proposed model. In the scheme of [7], inter-transmission times between fake transmissions are iid Exponentials with mean μ . The Anderson-Darling (A-D) goodness of fit test [15] is used to determine the time for transmitting real events without violating the exponential distribution of fake transmissions.³ Similarly, the A-D test is also used to implement the mean recovery algorithm. The authors of [7] used different statistical tests, such as the Kolmogorov-Smirnov (K-S) test [18], to show that their design satisfies event indistinguishability.

3. The Anderson-Darling goodness of fit test is a statistical test that, given a sequence of data samples and a desired degree of accuracy, determines whether the samples follow a certain probabilistic distribution with a certain parameter, within the specified degree of accuracy, or not.

4.3.1. Experimental Setup. For a reliable analysis of [7], we use the same parameters appeared in their paper. The inter-transmission times between fake message transmissions are set to be iid exponentials with mean 20 seconds. Real events arrive according to a Poisson Arrival process with mean 1/20. The two parameters of the A-D test are the significance level of the test and the allowed deviation from the mean which are set to 0.05 and 0.1, respectively.

The experiment was run for 10,000 independent trials. Each trial consists of two intervals, a real one I_R and a fake one I_F . Every trial starts with a “warmup” period, where 200 iid exponential random variables with mean 20 are drawn to constitute a backlog to be used in the A-D goodness of fit test. Then real events start arriving and they are transmitted according to procedure described earlier (please refer to [7] for detailed description and algorithms). Each real interval consists of 50 real events. After the 50th real event has been transmitted, the fake interval starts for the same amount of time the real interval lasted.

4.3.2. Simulation Results. After running the above experiment for 10,000 trials, and comparing the number of short-long patterns in fake and real intervals for each trial, the following results were found. Out of the 10,000 trials, real intervals have *more* short-long patterns than fake intervals in 6,818 trials; real intervals have *less* short-long patterns than fake intervals in 2,076 trials; and real intervals have the same number of short-long patterns as fake intervals in 1,106 trials.

4.3.3. Λ -anonymity Interpretation. Recall that, by equations (2) and (7), a short-long pattern is most likely to occur in real intervals than fake ones. Consequently, real intervals are likely to have more short-long patterns than fake intervals. Indeed, our simulation results coincide with equations (2) and (7).

Consider Game 1 for analyzing interval indistinguishability. Given two intervals I_0 and I_1 at which one of them is real and one is fake, let the adversary’s strategy for deciding which is which is as follows. Count the number of short-long patterns in each interval. If both intervals have the same number of short-long patterns, the adversary decides randomly. If one interval has more short-long patterns than the other, the adversary chooses it as the real interval. With this strategy, given the simulation results provided above, the adversary’s probability of correctly identifying real intervals, without resorting to complicated statistical tests, is 0.737. That is, the anonymity of the system is only $\Lambda = 0.526$.

5. New Direction to Improve Anonymity

So far, we have shown, in Section 4, that EI-based designs, although shown to provide location privacy under

existing models, do not provide high anonymity when analyzed under the proposed model. In this section, we suggest a new general method for designing transmission algorithms that can improve the anonymity of sensor networks. We will start by describing an overview of our approach. Then, we will provide a concrete example on how to apply the proposed approach on the EI-based scheme analyzed in the previous section and quantitatively compare the anonymity of the original design with the improved one. The example is not meant to be a proposed solution,⁴ it merely illustrates how the new direction can be applied to an existing EI-based design and quantifies the improvement in anonymity that can be achieved.

5.1. Overview

As can be seen from the analysis of the EI-based approach [7] in Section 4, inter-transmission times during fake intervals are iid's, while inter-transmission times during real intervals are neither independent nor identically distributed. This observation was the main factor behind increasing the adversary's chances in distinguishing between fake and real intervals in EI-based approaches. In theory, the only way to guarantee that a sequence of random variables is statistically indistinguishable from a given iid sequence is to generate it as an iid sequence with the same distribution. This implies that the only EI-based solution that guarantees absolute anonymity is the trivial solution of transmitting every real event in place of its successive scheduled fake message. As discussed earlier, however, the trivial solution does not minimize latency for arbitrarily distributed arrival of real events.

The notion of interval indistinguishability, apart from being the key point enabling the analysis of EI-based approaches, suggests a different approach for the design of anonymous sensor network systems. Observe that Definition 2 of interval indistinguishability does not impose any requirement, such as iid, on the distribution of inter-transmission times during fake intervals. That is, the inter-transmission times during fake intervals can have any arbitrary distribution. Therefore, designing fake intervals with the distribution that is easiest to emulate during real intervals is the most logical solution. In fact, since the arrival distribution of real events is generally not iid, it is only natural to design fake intervals with non iid inter-transmission times. This idea opens the door for more solutions as it gives more flexibility for system designers.

We suggest the following method for transforming EI-based designs into II-based to improve their anonymity. Instead of designing the transmission algorithm of real events based on a pre-fixed distribution for fake intervals,

4. The proposal of an efficient anonymous system based on our framework will be the focus of a future work.

the system can be designed as follows: given the desired algorithm for handling real events, fake intervals can be designed accordingly. That is, we suggest introducing the same correlation of inter-transmission times during real intervals to inter-transmission times during fake intervals. In what follows, we give a detailed example of how to apply this approach on the system analyzed in Section 4.3.

5.2. Concrete Example

Consider the same algorithm for real event transmission appeared in [7]. That is, when real events occur, their transmission time is computed as the minimum value that passes the A-D goodness of fit test. Furthermore, the transmission following a real event is delayed to adjust the ensemble mean. The fundamental problem here is that inter-transmission times in real intervals are correlated by design, and the example in Section 4.3 illustrates how this correlation can be exploited to reveal location information.

Therefore, as opposed to the scheme of [7], we design fake intervals to be as close as possible to real intervals. We suggest the generation of "dummy events" during fake intervals that are to be handled as if they are real events. That is, dummy events are generated independently from fake messages and, upon their arrival, their transmission times are determined according to the used statistical test. The purpose of this procedure is to introduce the same correlation of inter-transmission times during real intervals to the inter-transmission times during fake intervals.

However, recall that if the distribution of the arrival of real events is known, it is easy to design anonymous systems. Therefore, it is critical that the generation of the dummy events is independent of the distribution of real events. That is, the suggested approach must be doable without prior knowledge of the distribution of real events. The example below shows how the same tool used to design EI-based schemes, statistical goodness of fit tests, can be utilized to implement the suggested approach.

5.2.1. Setup. We adopted the same real interval transmission algorithm and parameters of [7] described in Section 4.3. That is, real events arrive according to a Poisson process with mean 1/20 and the inter-transmission times between fake messages are iid exponentials with mean 20 seconds. During fake intervals, fake messages are also scheduled as iid exponentials with mean 20 seconds.

To resemble real intervals, however, we generated *dummy events* according to iid Gaussian inter-arrival times with mean 10 seconds and a variance of 150. Note the distinction between fake messages and dummy events. Fake messages are the ones transmitted to hide the existence of real transmissions, while dummy events are the ones generated, during fake intervals only, to resemble the

Table 2. A quantitative comparison of the three schemes, the EI-based approach of [7], our II-based application of [7], and the trivial solution of sending real events instead of their successive scheduled fake transmissions. $I_R > I_F$ denotes more short-long patterns in real intervals, $I_R < I_F$ denotes more short-long patterns in fake intervals, while $I_R = I_F$ denotes equal short-long patterns in real and fake intervals. The simulation results are obtained from 10,000 independent trials.

	$I_R > I_F$	$I_R < I_F$	$I_R = I_F$	Anonymity of the system (Λ)
EI-based approach	6,818	2,076	1,106	0.526
Our II-based approach	4,566	4,272	1,162	0.971
Trivial solution	4,385	4,318	1,297	0.993

existence of real events. Furthermore, note that the inter-arrival distribution of dummy events is purposely different than the inter-arrival distribution of real events to count for the general case of unknown distribution of real events inter-arrival times.

Dummy events are handled as if they are real events. That is, in fake intervals, fake messages are transmitted according to iid exponential inter-transmission times and, upon the arrival of a dummy event, its transmission time is determined to satisfy the A-D goodness of fit test for a sequence of iid exponentials with mean 20 seconds.

Remark 1: As we mentioned earlier, this example is not meant to be a practical solution to the anonymity problem as it requires nodes to perform the A-D test even in the absence of real events. The example merely illustrates how one can transform the EI-based schemes of [7] into an II-based. In fact, the A-D test is proposed in the EI-based approach of [7] to make the overall transmission statistically similar to the iid distribution of pure fake transmissions. Since in our model fake intervals are not restricted to be iid, we believe that real and fake intervals can be similar without resorting to computationally cumbersome statistical tests.

Therefore, we do not provide thorough simulation analysis showing the effect of different distributions of dummy events. However, there is a supporting evidence suggesting that changing the distribution of dummy events will not have a considerable effect on anonymity of the system. To see this, recall that it has been shown in [7] that using statistical goodness of fit tests to handle real events transmissions make the overall transmission indistinguishable from the desired pre-specified distribution, regardless of the distribution of inter-arrival times of real events.

5.2.2. Simulation Results. After running the above experiment for 10,000 trials, and comparing the number of short-long patterns in fake and real intervals for each trial, the following results were found. Out of the 10,000 trials, real intervals have *more* short-long patterns than fake intervals in 4,566 trials, real intervals have *less* short-long patterns than fake intervals in 4,272 trials, and real intervals have the same number of short-long patterns as fake intervals in 1,162 trials.

To serve as a reference point for our anonymity comparison, we also simulated the trivial solution, where real events are transmitted instead of their successive scheduled fake messages. Out of the 10,000 trials, real intervals have *more* short-long patterns than fake intervals in 4,385 trials, real intervals have *less* short-long patterns than fake intervals in 4,318 trials, and real intervals have the same number of short-long patterns as fake intervals in 1,297 trials.

5.2.3. Λ -anonymity Interpretation. Consider Game 1 of analyzing interval indistinguishability. Given two intervals I_0 and I_1 at which one of them is real and one is fake, let the adversary’s strategy for deciding which is which be as described in Section 4.3.3. With this strategy, given the simulation results provided above, the anonymity of the improved II-based approach is $\Lambda = 0.971$, while it is $\Lambda = 0.993$ in the trivial solution. That is, an adversary basing her decision on the count of short-long patterns will be successful 51.5% of the time in the improved approach, with about 0.3% margin of error. Table 2 summarizes the performances of the three schemes.

Remark 2: Event indistinguishability is not included in the simulation since it follows from the A-D test. That is, transmitting dummy and real events according to the A-D test guarantees that inter-transmission times of each interval (fake or real) are statistically indistinguishable from the desired exponential distribution of fake messages only.

Observe the increased anonymity from 0.526 in the original EI-based approach of [7] to 0.971 in the improved II-based approach. This is obviously the desired behavior since it translates to lower probability of location detection. The way the improved II-based approach translates to higher anonymity is by increasing the number of short-long patterns in the fake interval, whereas the original EI-based approaches can only attempt to lower the confidence by decreasing the number of short-long patterns in the real interval, since fake intervals are fixed.

Remark 3: Observe that the notion of short-long patterns is merely a way of representing the class of correlation attacks to distinguish between real and fake intervals. Observe further that the improved approach does not

increase anonymity by directly increasing the number of short-long patterns in fake intervals. The anonymity is improved by trying to induce the same correlation between the inter-transmission times of real intervals to the inter-transmission times in fake intervals. This indirectly leads to more short-long patterns in fake intervals and, ultimately, to improved indistinguishability.

6. Conclusion and Future Work

In this paper, source anonymity in wireless sensor networks is addressed. We provided a statistical framework for modeling, analyzing, and evaluating anonymity in sensor networks. We introduced the notion of interval indistinguishability, proved that it implies the currently adopted model (event indistinguishability), and showed that it captures the source of information leakage that was not captured by event indistinguishability. Thus, the proposed anonymity model is stronger than existing models and allows for more rigorous anonymity analysis. We analyzed an EI-based approach, which was shown to provide anonymity under event indistinguishability, and quantified its information leakage when analyzed under our proposed model. Finally, we proposed a new direction for designing transmission algorithms that can improve source anonymity in sensor networks, applied our approach to an existing scheme, and quantified the improvement in anonymity that can be achieved.

Future extensions to this work include taking advantage of the key point that fake intervals are not restricted to have iid inter-transmission times to design an efficient system that satisfies the notion of interval indistinguishability, without resorting to computationally cumbersome statistical tests.

References

- [1] P. Kamat, Y. Zhang, W. Trappe, and C. Ozturk, "Enhancing Source-Location Privacy in Sensor Network Routing," *ICDCS 2005. The 25th IEEE International Conference on Distributed Computing Systems*.
- [2] C. Ozturk, Y. Zhang, and W. Trappe, "Source-location privacy in energy-constrained sensor network routing," in *Proceedings of the 2nd ACM workshop on Security of ad hoc and sensor networks*, 2004.
- [3] Y. Xi, L. Schwiebert, and W. Shi, "Preserving source location privacy in monitoring-based wireless sensor networks," in *IPDPS 2006. The 20th International Parallel and Distributed Processing Symposium*, 2006.
- [4] B. Hoh and M. Gruteser, "Protecting Location Privacy Through Path Confusion," in *SecureComm 2005. First International Conference on Security and Privacy for Emerging Areas in Communications Networks.*, 2005.
- [5] Y. Ouyang, Z. Le, G. Chen, J. Ford, F. Makedon, and U. Lowell, "Entrapping Adversaries for Source Protection in Sensor Networks," in *Proceedings of the 2006 IEEE International Symposium on World of Wireless, Mobile and Multimedia Networks*, 2006.
- [6] K. Mehta, D. Liu, and M. Wright, "Location Privacy in Sensor Networks Against a Global Eavesdropper," in *ICNP 2007. IEEE International Conference on Network Protocols.*, 2007.
- [7] M. Shao, Y. Yang, S. Zhu, and G. Cao, "Towards Statistically Strong Source Anonymity for Sensor Networks," *INFOCOM 2008. The 27th IEEE Conference on Computer Communications.*, 2008.
- [8] Y. Yang, M. Shao, S. Zhu, B. Urgaonkar, and G. Cao, "Towards event source unobservability with minimum network traffic in sensor networks," in *Proceedings of the first ACM conference on Wireless network security*, 2008.
- [9] N. Li, N. Zhang, S. Das, and B. Thuraisingham, "Privacy preservation in wireless sensor networks: A state-of-the-art survey," *Ad Hoc Networks*, 2009.
- [10] H. Wang, B. Sheng, and Q. Li, "Privacy-aware routing in sensor networks," *Computer Networks*, 2009.
- [11] M. Shao, W. Hu, S. Zhu, G. Cao, S. Krishnamurthy, and T. La Porta, "Cross-layer Enhanced Source Location Privacy in Sensor Networks," *SECON 2009. Sixth Annual IEEE Communications Society Conference on Sensor, Mesh and Ad Hoc Communications and Networks.*, 2009.
- [12] D. Chaum, "Untraceable electronic mail, return addresses, and digital pseudonyms," *Communications of the ACM*, 1981.
- [13] M. Reed, P. Syverson, and D. Goldschlag, "Anonymous connections and onion routing," *IEEE Journal on Selected Areas in Communications*, 1998.
- [14] M. Gruteser and D. Grunwald, "Anonymous Usage of Location-Based Services Through Spatial and Temporal Cloaking," in *Proceedings of the 1st international conference on Mobile systems, applications and services*, 2003.
- [15] M. Stephens, "EDF statistics for goodness of fit and some comparisons," *Journal of the American Statistical Association*, 1974.
- [16] M. Shao, S. Zhu, W. Zhang, and G. Cao, "pDCS: Security and Privacy Support for Data-Centric Sensor Networks," *INFOCOM 2007. The 26th IEEE International Conference on Computer Communications.*, 2007.
- [17] S. Goldwasser and S. Micali, "Probabilistic encryption," *Journal of Computer and System Sciences*, vol. 28, no. 2, pp. 270–299, 1984.
- [18] F. Massey Jr, "The Kolmogorov-Smirnov test for goodness of fit," *Journal of the American Statistical Association*, vol. 46, no. 253, pp. 68–78, 1951.