# How to Detect Cloned Tags in a Reliable Way from Incomplete RFID Traces

Mikko Lehtonen*, Florian Michahelles*, and Elgar Fleisch*†

*Information Management, ETH Zürich, 8092 Zürich, Switzerland, {mikkol,fmichahelles}@ethz.ch
†University of St.Gallen, Institute of Technology Management, 9000 St. Gallen, Switzerland, elgar.fleisch@unisg.ch

*Abstract*—**Cloning of RFID tags may lead to considerable financial losses and worse reputation in many commercial applications, while being attractive for adversaries. One way to address tag cloning is to use the visibility that RFID traces provide to detect cloned tags as soon as they enter the system. However, RFID traces always represent historic events without giving certainty where the traced objects currently really are. Furthermore, imperfect read rates can lead to missing reads. As a result, the visibility is not always perfect, which makes detection of cloned tags harder and less reliable. This paper presents a series of probabilistic techniques to enable reliable detection of cloned tags in cases where the visibility is incomplete. Our hypothesis is that the events generated by cloned tags cause rare or abnormal events that can be detected when the process that generates the legitimate events is understood. The presented techniques are studied in a comprehensive simulation study of a real-world pharmaceutical supply chain. Our findings suggest that reliable detection of cloned tags is possible if missing reads are addressed and the supply chain is precisely modeled.**

## I. INTRODUCTION

Radio Frequency Identification (RFID) systems are used to identify physical objects in many applications where cloning of tags could lead to considerable financial losses and harm. Examples of such applications include access control [1], ticketing [2], payment [3], and anti-counterfeiting in supply chains [4]. The conventional approach to secure RFID systems against tag cloning is to use cryptographic tags that enable tag authentication and make tag cloning considerable harder [5]. The less conventional approach to address tag cloning is to use the visibility that RFID provides to detect cloned tags.

Visibility enables location-based product authentication that is a technical anti-counterfeiting measure that brand-owners can use in their fight against product counterfeiting. Today, it is relatively easy to produce visual copies of various kinds of trademarked and branded products. Selling these counterfeit goods as originals to the licit supply chains can result into high illegal profits for the illicit players. Overall, counterfeiting constitutes a large problem for legally run businesses and governments; for example, the European Customs alone seize up to a hundred million counterfeit articles every year [6].

An important problem behind location-based product authentication is that RFID traces always represent historic events without giving certainty where the traced objects currently really are. We define an RFID trace as the set of events relating to one tagged object that can be retrieved from a tracing system, such as the EPC network [22]. Furthermore, the existing RFID systems are still somewhat prone to read errors. As a result, the visibility that RFID systems provide is not always perfect, which makes detection of cloned tags harder and less reliable. We propose to address the problem of incomplete traces with probabilistic analysis techniques whose goal is to automatically detect cloned tags in a reliable way from a large amount of data. We make use of the fact that events generated by cloned tags appear in traces of genuine products. Our hypothesis is that the events generated by cloned tags cause rare or abnormal events that can be detected as improbable transitions in the supply chain when the process that generates the legitimate events is understood and modeled. If this hypothesis is true, it implies that supply chains can be protected from cloned tags with detective techniques, without the need for cryptographic tags.

This paper is organized as follows. Section II reviews related work about uncertainty in RFID traces and techniques to address tag cloning attacks. In Section III we present the characteristics of RFID track and trace data, including an analysis of causes of missing reads. Sections IV and V present the probabilistic reasoning behind location-based product authentication and the proposed solution method. In Section VI we evaluate the proposed methods with a simulation study of a real-world pharmaceutical supply chain. We finish with discussion and conclusions.

## II. RELATED WORK

Various causes of uncertainty in RFID traces have been discussed in the scientific community. Derakhshan *et al.* [7] identified inaccuracy, i.e. missing reads (false negatives) due to imperfect read rates, as one of the primary factors limiting the widespread adoption of RFID technology. The authors stated that the observed real-world read rates are often in the 60-70% range [8]. However, presenting general numbers for read rate is misleading since read rate depends on a multitude of case-specific factors. Case studies in the aerospace industry [9] show that in addition to missing reads, also processing delays and problems with aggregation information accuracy can decrease the tracking information quality.

Examples from real life show that missing reads can be eliminated in real-world RFID implementations by engineering solutions in the physical layer. Folcke [10] described a commercial "smart cabinet" application where medical devices (25% of which contained a metal cover) are automatically identified. After the first deployment, the error rate in identification was smaller than $1.5 \cdot 10^{-5}$ (no errors in 65,000

reads). This high level of reliability was achieved by choosing the best frequency for that particular application (125 kHz), by positioning and adjusting the reader antennas, and by choosing the best tag position on the products.

Researchers in the University of Arkansas measured the read rates of RFID tags in Electronic Article Surveillance (EAS) [11]. Overall, the tested UHF tags and readers performed very well and compared to existing EAS systems, especially for reading single tags, and displayed a general insensitivity to tag orientation. When scanning 50 tags passing through a gate, read rates of 95% and more were observed. Brusey *et al.* [12] studied the problems of uncertainty in the true locations of tags in industrial robotic control application and smart medicine cabinet application. The authors found that filtering tag reads over time was effective in reducing both false positive and false negative reads. Jeffery *et al.* [14] proposed the first adaptive smoothing filter for RFID data cleaning generated by the readers. In another work, the same authors proposed a data cleaning process for sensor/RFID data to solve missed reads when scanning multiple products on a shelf [13].

Khoussainova *et al.* [15] studied how to address uncertainty in the application layer. As a solution to difficulties in determining which high-level event has occurred based on the low-level reads, the authors proposed a probabilistic model that divides uncertain events into multiple high-level events with given probabilities. Kelepouris *et al.* [16] studied the quality of tracking information in supply chains and discussed the problem that RFID data does not tell the products' current locations without uncertainty. The authors introduced the first quantitative metric of the quality of product location information based on expected utility axiom.

The research community addresses authentication of RFID-tagged products primarily by trying to make tag cloning hard using cryptographic tag authentication protocols [5]. The fundamental difficulties of this research revolve around the trade-offs between tag cost, level of security, and performance in terms of reading speed and distance. It is not sure if the cryptographic solutions will be inexpensive and usable enough for deployment as barcode replacing RFID tags. Furthermore, solving key distribution and management in a scalable and secure way is challenging for large RFID systems.

In addition to cryptographic approaches, various ways to detect counterfeit products based on track and trace data have been discussed in the literature. Electronic pedigree(e-pedigree), where buyers and sellers of the product append the product's history document with an event that they sign, is probably the best-known measure [17]. The major limitation of e-pedigree is that it does not provide a reliable way of detecting copied e-pedigree documents. Juels [5] noted that serial level identification alone without secure verification of the identities can be a powerful anti counterfeiting tool. Koh *et al.* [18] made use of this assumption to secure pharmaceutical supply chains by proposing an authentication server that publishes a list of genuine products' ID numbers. Takagari *et al.* [19] proposed some early ideas how to check the validity of serialized ID numbers. Staake *et al.* [4] were among the
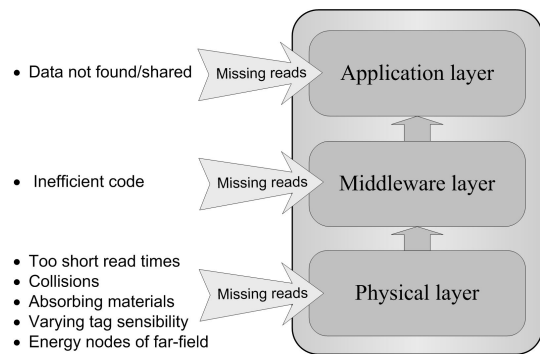


Fig. 1. Sources of missing reads in RFID systems

first to discuss the potential of track and trace based product authentication and they point out some problems that occur when the back-end no longer knows where the genuine subject is. Mirowski and Hartnett [1] developed a system that detects cloned RFID tags, or other changes in tag ownership, in an access control application with intrusion detection methods.

Lehtonen *et al.* [20] used machine learning techniques to automatically detect cloned tags from incomplete location data. The authors applied time-delayed Markov chains and hidden Markov models (HMM) to classify traces of products as either clean of cloned tags or corrupted by cloned tags. Their findings suggested that cloned tags can be best detected by searching for single unlikely transitions from track and trace data. This paper uses a similar method, but it is applied to single events instead of complete traces. Furthermore, this work includes a method to detect missing read events and presents a detailed evaluation of the authentication method based on a real-world supply chain.

### III. CHARACTERISTICS OF RFID TRACES

This section summarizes the characteristics of RFID traces from those parts that are relevant to anti-counterfeiting. In general, the semantic attributes of events follow a "*What? When? Where? Why?*" -concept. For an anti-counterfeiting application, the most relevant information of events is captured by the time and location location attributes. The event time is simply the time when the event occurred and was captured. The event data of EPCIS 1.0.1 specification [21] defines physical and logical reader attributes, a read point attribute, and a business location attribute. Among these attributes, the business location is the most suitable for location tracking since it defines the discrete and unambiguous location where the object is after the event.

Moreover, RFID events do not always tell the object's current location, but where the object has been observed. As a result, track and trace data includes uncertainty about the object's current location. Missing reads contribute to this *location uncertainty* and thus decrease the usefulness and value of track and trace data. The possible causes of missing reads are analyzed below and summarized in Fig. 1.

Most causes of missing reads can be traced down to the

physical layer. Most importantly, these include i) missing reads due to too short reading times (the time the tag is in the reader's field), ii) collisions in the air interface that collision detection protocol does not catch, and iii) conductive materials that absorb radio waves. In addition, when tags are read using the far-field, the tag might be in a node where the field strength is close to zero and thus the tag will not be read. Furthermore, due to normal variance in tag manufacturing processes, chips and antenna connections have varying impedances which results into variance in tag read ranges. It is important to note that bit errors in the data that is read from the tags do not constitute a source of uncertain event attribute values because of error detection coding.

In some cases slow or otherwise not optimal code in the middleware can result into missing reads even when an antenna has interrogated the tag. The role of the middleware is to coordinate multiple readers that occupy the same physical space and to transform raw tag reads into streams of high-level events for example by filtering, aggregating, and counting them [21]. In a typical setting, middleware listens to all the antennas of a reader device inside a loop. If a tag is present in one antenna's field whilst the middleware is listening to other antennas (e.g. reading other tags), the low-level event might not be captured by the middleware.

Also problems in the application layer can lead to missing events. For instance, Discovery Services [22] might not be able to locate events relating to a product with a 100% reliability.

In addition, RFID traces can be plagued by so called *phantom reads* where a reader reports a tag that was was not in its field or did not exist. Phantom reads, however, are not considered in the remainder of this paper.

## IV. LOCATION-BASED PRODUCT AUTHENTICATION

Tracking and tracing enables location-based authentication [20]. The underlying assumptions are that all genuine products have a unique ID number and there exists a way to find out whether a unique ID number is valid or not. This scheme is not yet secure because an adversary could clone a tag. The location-based authentication system secures this scheme by detecting the cloned tags based on their locations.

Detection of cloned products from the track and trace data is straightforward if the current locations of the products are precisely known; for instance, if the track and trace data tells that the product is currently in Switzerland at the same time when a product with the same ID is scanned in Japan, the system can conclude that it is probable that the product in Japan has a cloned tag. However, when the track and trace data says that the product was observed in Switzerland one week ago but it does not tell its current location, authentication becomes harder and false alarms become possible.

We build an automatic location-based authentication system by evaluating *transition probabilities* between the events. A transition probability stands for the probability that a genuine product makes the transition defined by two events. If the transition probability is high, the latter event is likely to be generated by a genuine tag and vice versa. When we denote event $i$

as $E^i$, the transition probability $P_{tr}$ from $E^i$ to a consecutive event $E^{i+1}$ can be presented as $P(E^{i+1}|E^i, E^{i-1}, ..., E^1)$ (cf. Fig. 2). As a result, the authentication rule can be formalized as follows. *Event $E^i$ is generated by a genuine product if*:

$$P_{tr} = P(E^{i+1}|E^i, E^{i-1}, ..., E^1) > \epsilon \qquad (1)$$

The transition probability of the first event ($i = 1$) in a product's trace can be estimated by introducing a so called "zero-event". Like this, the transition probability of the first event is given by $P(E^1|E^0)$. By limiting the locations where this probability is non-zero, the system defines a limited secure environment where new products are allowed to occur (e.g. a manufacturer's packaging line).
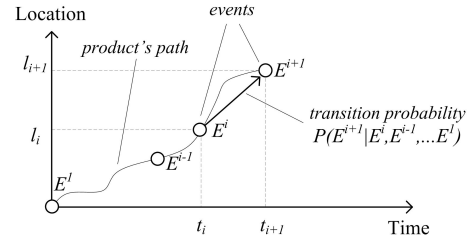


Fig. 2. Events and transition probabilities

The location-based authentication problem can be now solved by building a selective classifier that yields high transition probabilities for events generated by genuine tags and low transition probabilities for events generated by cloned tags. Using the Bayes' rule, the transition probability can be further turned into an *a posteriori* probability that an event is generated by a genuine tag. We denote the transition probability generated by an event by the random variable $X$. The probability that an event is generated by a genuine product, $P(ge)$, given that the transition probability the event generated is smaller than $x$, can be formulated as follows:

$$P(ge|X < x) = \frac{P(ge, X < x)}{P(X < x)} \qquad (2)$$

$$= \frac{P(ge) \cdot P(X < x|ge)}{P(X < x)} \qquad (3)$$

All terms in the last expression can be estimated from data that contains known counterfeit products.

## V. PROBABILISTIC SOLUTION METHOD

In this section, we present our probabilistic solution method that instantiates the authentication approach outlined above. The data processing steps of our solution are following:

1) Train the supply chain model with training data,
2) Filter the testing data set to find missing reads,
3) Evaluate $P_{tr}$ for all events in the filtered data, and
4) Raise an alarm if $P_{tr}$ is below a threshold.

We formulate the generic transition probability (Equation 1) into a more useful form. Two attributes are enough to give a semantically rich presentation of RFID events, namely event

time ($t$) and the discrete business location ($l$) [21]. We start with the first order Markov assumption which says that the state of the system is fully described by the last event, or:

$$P(E^{i+1}|E^i, E^{i-1}, ..., E^1) = P(E^{i+1}|E^i) \qquad (4)$$

This assumption discards path dependency of business locations. By assuming that time and location of new events are mutually independent random variables, and that locations of new events do not depend on time of the preceding events, we can express the transition probability as follows:

$$
\begin{aligned}
P(E^{i+1}|E^i) &= P(l_{i+1}, t_{i+1}|l_i, t_i) \qquad (5)\\
&= P(l_{i+1}|l_i, t_i) \cdot P(t_{i+1}|l_i, t_i)\\
&= P(l_{i+1}|l_i) \cdot P(t_{i+1}|l_i, t_i)\\
&= P_{i,i+1} \cdot P(\Delta T_i = t_{i+1} - t_i)
\end{aligned}
$$

### A. Stochastic supply chain model

To evaluate the two terms in the last expression of Equation 5, we model the process how track and trace events are generated in a supply chain. We model the supply chain as nodes and lines and build a Stochastic Supply Chain Model (SSCM) that has $N+1$ distinct states, $S_0, S_1, S_2, ..., S_N$. The relation between states in the model and the observed events is following: every time a product enters a state in the model, it generates a track and trace event in the real life. In other words, a state in the model corresponds to a reader device. The zero-state, $S_0$, represents the "state of non-existence" where all tagged products are before they are created in the real world, and exceptionally it does not have corresponding events or business location in the real life. All other states in the model correspond to discrete business locations of the real-world supply chain network where tagged products are read. Parameters of the model define how products move from one discrete business location to another.

In the common case, after entering a state, the product stays there during a finite number of steps. This corresponds to a normal observation event. The time before the product generates a new event, called the waiting time, is given by a probability density function (PDF) that is specific to each state. For state $i$, $1 \leq i \leq N$, this PDF is denoted as $p(\Delta T_i)$. The actual distribution is not constrained by the model and it can be e.g. uniform or Gaussian. After time $\Delta T$ from entering a state, the product enters a new state according to the state transition probabilities. The first event in a product's trace is generated when the product leaves the zero-state $S_0$. After that, the product continues to move in the model through normal states as described above until it reaches an end-state. There are no routes that leave an end-state and thus the waiting time in an end-state can be regarded as infinite.

The state transition probabilities are time independent and denoted as $P_{ij} = P(S_i|S_j) \geq 0$, $i, j \geq 0$. State transition probabilities from a state to itself ($P_{ii}$) are possible and they correspond the real-life situation where a product's trace has two consecutive reads from the same single business location.
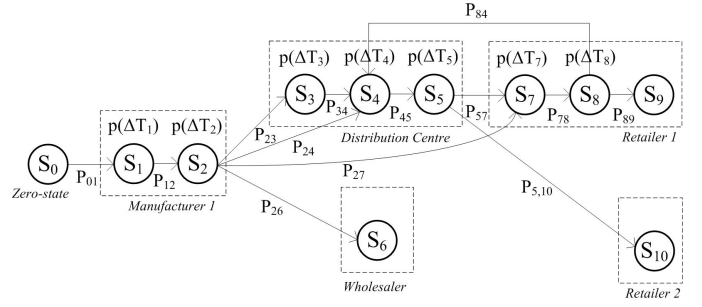


Fig. 3. Illustration of our Stochastic Supply Chain Model (SSCM)

Each physical location in the supply chain is represented in the SSCM by three states corresponding to receiving, internal, and shipping operations. The SSCM is trained from RFID traces and therefore only locations where products are scanned are present in the SSCM. The resulting model is flexible and intuitive and it has enough degrees of freedom to capture the essential statistics of how single products flow in supply chain networks. The SSCM is exemplified in Fig. 3. This imaginary supply chain illustrates different real-world problems in location-based product authentication: missing reads at reader in business location $S_3$ (results into a "ghost route" $P_{24}$, cf. subsection V-B), a wholesaler and a retailer that do not share trace data beyond receiving notifications ($S_6$ and $S_{10}$, respectively), and reverse logistics ($P_{84}$).

If the model would use state transition probabilities from a state to itself to define the time a product stays in a state instead of the waiting time PDFs, the model would be a time-independent first-order discrete time Markov chain (DTMC). However, we have opted for defining the waiting time distribution because it allows for flexible modeling of the supply chain's time dynamics (i.e. in DTMC the waiting time distribution is fixed while in SSCM it can have any form).

### B. Filtering traces to detect missing reads

The SSCM can be used to detect missing reads (cf. Section III) in RFID traces. Missing reads can trigger unwanted false alarms in the clone detection system. Reader devices that have a below 100% read rate create "ghost routes" that are observed as small transition probabilities that do not correspond to real-world transitions (cf. Fig. 4). Our filtering algorithm tries to detect when a product is moving along such a "ghost route" as evidence of a missing read event.

We explain how the filtering detects missing read events by referring to the example in Fig. 4. When a transition probability is low (from A to C), the filtering algorithm can search for a more probable alternative route that is obtained by including a new read event between the existing events. If the probability of the new route (from A to B to C) is higher than a threshold, the new event (in B) is added to the trace.

The number of missing consecutive read events that the filter can add is called the order of the filter. In this paper we study 1st and 2nd order filters. Filters of all orders can be described by three parameters: i) maximum transition
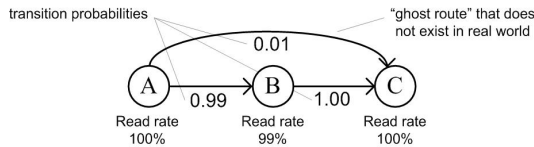
Fig. 4. Though all products flow from A to B to C, the 99% read rate in B creates a "ghost route" A to C. The filtering algorithm tries to detect when a product moves along the "ghost route" to add the missing read event in B.

probability (threshold) between the existing events, ii) minimum time difference (threshold) between the existing events, and iii) minimum geometric mean (threshold) of transition probabilities of the new route. The first two parameters define when the filter is allowed to add missing reads between existing events and the third parameter limits the addition of new routes that are too unlikely. The values of these parameters can be defined empirically.

### C. Location-based authentication

For $E^i$, $i > 1$, SSCM enables evaluation of a location transition probability ($P_{i-1,i}$) and a time transition probability ($P(\Delta T_{i-1} = t_i - t_{i-1})$). We denote these methods as $SSCM_L$ and $SSCM_T$, respectively, and we compare their performance in a simulation study. For the first event in a trace, $E^1$, only the location transition probability is defined. Now the authentication rule from Equation 1 can be rewritten in two new ways. *Event $E^i$ is generated by a genuine product if*:

$$SSCM_L: P_{i-1,i} > \epsilon \qquad (6)$$

$$SSCM_T: P(\Delta T_{i-1} = t_i - t_{i-1}) > \epsilon \qquad (7)$$

The value of the threshold $\epsilon$ defines the trade-off between the ratio of event of cloned tags that are detected (hit rate) and the ratio of events of genuine products classified as generated by cloned tags (false alarm rate). The value of $\epsilon$ can be optimized only by setting a cost for false alarms and a value for hits. In practice, minimization of false alarms might be wanted and hence $\epsilon$ can be set to the smallest transition probability of genuine products within the training data. In general, the threshold $\epsilon$ has different values in Equations 6 and 7.

We believe that an optimal location-based authentication system should somehow combine the location and time transition probabilities presented in Equations 6 and 7. Finding an suitable way to combine these probabilities is, however, out of the scope of this paper.

## VI. SIMULATION STUDY

We evaluate the proposed methods with a simulation study of a real-world pharmaceutical supply chain. The goal of this study is to evaluate how cloned tags that can be distinguished from the corresponding genuine tags in the presence of missing reads and a limited amount of training data. Cloned tags that appear before the corresponding genuine products are manufactured or after they are consumed are not considered because they can be detected with simple rules.

We measure the **hit rate**, i.e. how often events created by cloned tags are detected (system raises an alarm), versus the **false alarm rate**, i.e. how often alarms are triggered by events of genuine tags. The resulting trade-off is presented as a Receiver Operating Characteristics (**ROC**) curve that characterizes the selectivity of a classifier. In a real-world anti-counterfeiting application, only very small false alarm rates can be tolerated because the number of read events that the genuine products generate is very high.

Only the first events generated by the cloned tags are considered in the results. The reason is that the simulated supply chain handles both counterfeit and genuine products in an identical way, so the further events generated by cloned tags have identical statistics than events of genuine products. Thus the results indicate how reliably the cloned tags can be detected as soon as they enter the supply chain. In addition, those events of genuine products that are directly preceded by events from cloned tags are neglected from the results.
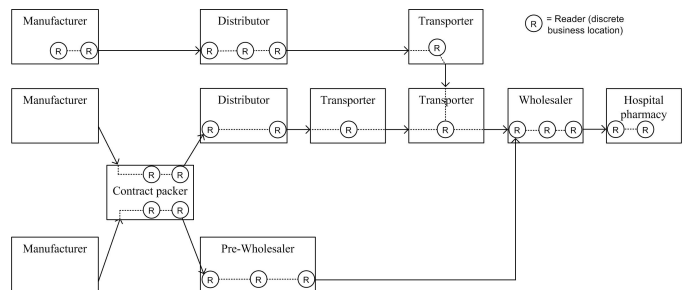


Fig. 5. The simulated real-world supply chain (R denotes a reader device)

### A. Simulated real-world pharmaceutical supply chain

The real-world pharmaceutical supply chain under study involves nine different organizations in the UK and Holland, including three manufacturers, a contract packer, distributors, a pre-wholesaler, and a wholesaler that supplies a hospital pharmacy in a major London hospital [23]. The products that flow through this supply chain are equipped with printed Data Matrix codes that store serialized ID numbers. Single packs are aggregated into cases and pallets that have both RFID tags and Data Matrix codes. The pallets are scanned in 20 read stations in different supply chain locations to generate track and trace events. The average lead time from production to hospital is about 40 days, varying between approximately one week and two months. The supply chain is illustrated in Fig. 5.

In the studied supply chain, traces of products begin either at the manufacturer's production line or at the contract packer's packaging line. Products are shipped to the wholesaler in pallets and the wholesaler uses a "pick, scan, and drop process" to fill boxes that fulfill the pharmacy's orders. The wholesaler delivers products to the hospital pharmacy 2-6 times a day according to orders. The last event in a product's trace occurs when it is scanned in to the hospital pharmacy's inventory, after which the products are identified based on the non-serialized EAN-13 bar codes.

We have built a model of the described supply chain in our own supply chain simulator. The simulator works with three-hour-long discrete time steps. The model is built based on documentation [23] and interviews and it has been validated with direct feedback and example track and trace data. In the simulator, each supply chain node is presented by three different locations corresponding to business steps of receiving, internal processes, and shipping. The time how long an object spends in these locations is given by a uniform distribution. If the product enters a location where there is a reader device, and no read error occurs, a track and trace event is generated. The transitions between the supply chain nodes are determined by transition probabilities. The transition times between the nodes are deterministic and estimated from the distances and transport methods (ship or truck).

The times that logistic units spend in different locations could not be accurate modeled since the real lead time distributions were not precisely known. However, more accurate modeling of the real-world lead times is not likely to affect the results. In addition, because we evaluate the transition probabilities without taking into account correlations among different products' traces, the simulator treats all logistic units as independent from each other, which means that for example aggregation events are not modeled.

### B. Set-up

One simulator run generates and analyzes one example set of RFID traces. In each run, all three manufacturers produce 500 tagged products per day during days 1 to 7. This creates 10,500 genuine products and more than 130,000 possible read events. During days 8 to 35, 8 counterfeit products are injected into randomly chosen non-manufacturer supply chain locations per day, constituting a total of 224 counterfeit products (resulting into a 2% counterfeit market share, a high but possible value for seriously infiltrated markets). The counterfeit products have ID numbers of randomly chosen genuine products so the events they generate appear in traces of 224 different genuine products. The simulation stops after 60 days. In some rare cases a counterfeit and a genuine product with the same ID are both scanned during the same time step. These cases are not considered in the results.

The results are calculated from the average ROC curves of 10 Monte Carlo iterations (i.e. simulator runs). Each iteration yields a number of discrete points in the ROC curve and a continuous curve is drawn by interpolating. The SSCM is trained in each iteration from the training data set and the waiting time distributions in the SSCM are uniform distributions between the smallest and biggest observed waiting times in that business location. The following tests are performed:

- **Test 1**: The performance of filtering algorithm in finding missing reads from trace data without cloned products with read rates 99.9%, 99.0%, 95%, and 90%, with training data size of 1000 traces.
- **Test 2**: The performance of $SSCM_L$ and $SSCM_T$ with read rates 99.9%, 99.0%, 95%, and 90%, with training data size of 300 traces.

- **Test 3**: The performance of $SSCM_L$ with with training data size 1000, 300, 100, and 50 traces, and read rates 99.9% and 99%.
- **Test 4**: The performance of filtering and $SSCM_L$ with 99% read rate and with training data size of 300 traces.

### C. Results

Results of Test 1 show that our filtering algorithm (subsection V-B) is able to detect up to 86% of missing read events, depending on the read rate and the filter order (cf. Table I). In practice it means, for example, that effective read rate can be increased from 99.0% to 99.84%. Second order filter is able to detect more missing reads than the first order filter when the read rate decreases because of the greater number of consecutive read errors. Moreover, the filter parameters were defined empirically, which leaves room for optimization.

Results of Test 2 show that that the location-based $SSCM_L$ is much more reliable in detecting cloned tags than the time-based $SSCM_T$ (Fig. 6). Overall, $SSCM_L$ provides reliable detection results, though the hit rates at the zero false alarm rate are less than 30%. Analysis of false alarms of $SSCM_L$ reveals that in cases when the cloned tag is injected into the location where the genuine product is expected, the cloned tag was not detected (miss) and the genuine product generated a false alarm. The tested $SSCM_T$ method is very prone to false alarms and thus it is not suitable in the studied clone detection application, but the form of the ROC curve confirms that also the transition times carry information that distinguishes events generated by cloned tags from normal events. The results of Test 2 also confirm that missing reads decrease the performance of the studied clone detection methods.

Results from Test 3 show that increasing the amount of training data improves the reliability of $SSCM_L$ in the presence of missing reads (Fig. 7). When the number of missing reads is small, a small amount of training data is enough for accurate modeling of the underlying supply chain. When the number of missing reads increases, more and more "ghost routes" (cf. Fig. 4) appear and more training data is needed to capture them. This indicates that precise modeling of the supply chain contributes to reliable detection of cloned tags.

Results from Test 4 show that our filtering algorithm decreases the number of false alarms caused by missing reads, increasing the hit rate at zero false alarm rate from zero to ca. 80% (Fig. 8). Analysis of misses reveals that in some rare cases the filter adds an event before the first event of a cloned tag, causing the miss. However, the overall effect of filtering is clearly positive. The posterior distribution in Fig. 8

TABLE I
NUMBER OF MISSED READ EVENTS WITH DIFFERENT FILTERS

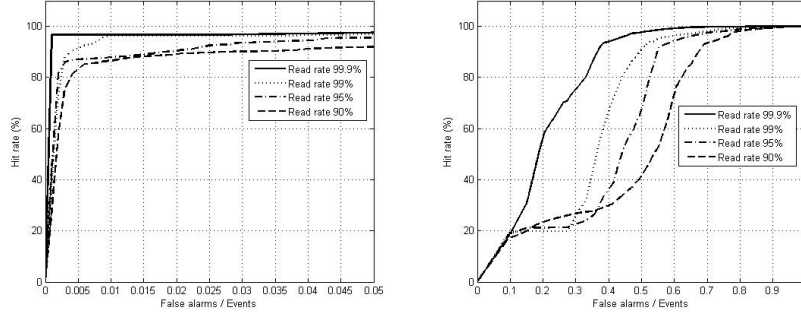| Read rate | No filter | 1. Order | 2. Order |
|---|---|---|---|
| 99.9% | 160 (100%) | 23 (14%) | 23 (14%) |
| 99.0% | 1392 (100%) | 246 (18%) | 228 (16%) |
| 95.0% | 6875 (100%) | 1541 (22%) | 1207 (18%) |
| 90.0% | 13920 (100%) | 4262 (30%) | 2821 (20%) |

Fig. 6. Results of Test 2: ROC curves for $SSCM_L$ (left) and for $SSCM_T$ (right) based clone detection. The curves show that $SSCM_L$ is much more reliable than $SSCM_T$ in detecting cloned tags, and that missing reads decrease the performance of both these methods. (note the different scales in x-axis)
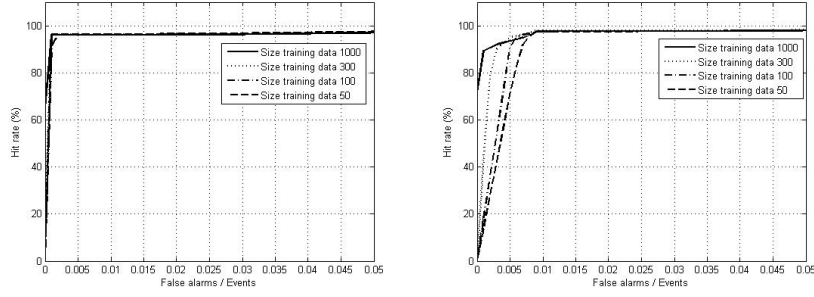


Fig. 7. Results of Test 3: ROC curves for $SSCM_L$ with 99.9% (left) and 99% (right) read rates. The curves show that increasing the amount of training data (more accurate modeling of the supply chain) is important for reliable detection of cloned tags as the number of missing reads increases.
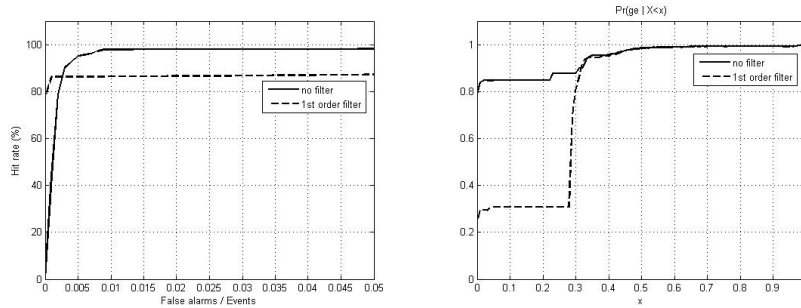


Fig. 8. Results of Test 4: ROC curves (left) and posterior distributions (right) of non-filtered and filtered traces for $SSCM_L$ with 99% read rate. The curves show that our filtering algorithm that detects missing reads can provide a dramatic increase to the hit rate with small false alarm rates.

proofs this by showing that the filtering algorithm increases the probability that an alarm is generated by a counterfeit product by about 50% in small false alarm rates.

## VII. DISCUSSION

Compared to other published results, the achieved above-95% hit rates at below-0.2% false alarm rates (cf. Fig. 6) indicate reliable detection of cloned tags. In a somewhat similar study, a 47% hit rate at a 1% false alarm rate was achieved by simpler supply chain modeling [20]. In an RFID-based access control system, cloned tags were detected with hit rates of 76%-46% at false alarm rates of 8.4%-2.5% [1].

The results of our simulation study confirm that the majority of cloned tags appear as abnormal events in RFID traces as soon as the tags enter the supply chain. This means that anomaly-based intrusion detection system techniques that are widely used to secure IT systems can be applied to detecting counterfeit products from track and trace data. Also missing reads that are common in today's RFID systems cause abnormal events and thus create false alarms, but they can be mitigated by our filtering algorithm that is able to detect up to 84% of missing reads. Moreover, our study shows that accurate modeling of the underlying supply chain contributes to reliable detection of cloned tags. The time transition probabilities did

not perform well in clone detection, but we still believe that event times include information that an optimal location-based authentication system should make use of.

The training data set needs to have an adequate quality by including all allowed transitions and no cloned tags. Instead of training the SSCM, a supply chain manager could alternatively set up the SSCM manually by selecting all allowed transitions and estimating the time distributions.

The concept of location-based authentication is not without limitations. If two products with the same ID are in the same location, a location-based authentication system cannot conclude which product is the genuine one. In addition, the system can generate false alarms that end-users need to deal with. Despite these shortcomings, the presented method presents a major complication to counterfeiters who want to inject counterfeit products into a licit supply chain. Most importantly, this countermeasure is based on processing of track and trace data, which does not increase the tag price and the tag reading time.

## VIII. CONCLUSIONS AND FUTURE WORK

In this paper, we present probabilistic techniques to detect cloned tags from RFID traces. The presented techniques enable detection of counterfeit products in supply chains where single products are traced. The results of our simulation study of a real-world pharmaceutical supply chain confirm that only in very exceptional cases cloned tags do not create unexpected events that can be detected. This finding implies that detection-based security measures have a very big potential to reliably detect cloned tags in well predictable processes, for example in a supply chain. Furthermore, we present a high-level event filtering technique to detect missing reads that constitute the biggest cause of false alarms in our clone detection application. Overall, the presented methods provide a considerable level of protection against serialized counterfeit products that enter a supply chain, without the need for cryptographic tags.

Future work towards an optimal track and trace based authentication system will investigate ways to combine the location transition probabilities with event times that carry complementary information. In addition, reliability of the system can potentially be further enhanced by taking into account correlations between different products' traces as well as information in related IT systems, such as advanced shipping notices.

## ACKNOWLEDGMENT

## REFERENCES

[1] Mirowski, L. T., Hartnett, J.: Deckard: a system to detect change of RFID tag ownership. Int. J. Comput. Sci. and Netw. Secur. 7(7), 89-98 (2007)
[2] Swedberg, C.: RFID Drives Highway Traffic Reports. RFID Journal. http://www.rfidjournal.com/article/view/1243/1/1 (2004). Accessed 3 December 2008
[3] Texas Instruments: ExxonMobil Speedpass. http://www.ti.com/rfid/shtml/apps-contactless-speedpass.shtml (2008). Accessed 3 December 2008
[4] Staake, T., Thiesse, F., Fleisch, E.: Extending the EPC network – the potential of RFID in anti-counterfeiting. In: Proceedings of Symposium on Applied Computing, ACM, 1607-1612. New York (2005)
[5] Juels, A.: RFID security and privacy: A research survey. IEEE J. Sel. Areas Commun. 26(2), 381-894 (2006)
[6] European Commission: Statistics recorded at the external borders of the EU. http://ec.europa.eu/taxation_customs/customs/customs_ controls/counterfeit_piracy/statistics/index_en.htm (2008). Accessed 30 August 2008
[7] Derakhshan, R., Orlowska, M., Li, X.: RFID data management: Challenges and opportunities. In: Proceedings of IEEE International Conference on RFID, 26-28. Grapevine, Texas (2007)
[8] Jeffery, S., Garofalakis, M., Franklin, M.: Adaptive cleaning for RFID data streams. In: Proceedings of the 32nd International Conference on Very Large Databases (VLDB), 163-174. Korea (2006)
[9] Kelepouris, T., Da Silva, S., McFarlane, D.: Automatic ID systems: enablers for track and trace performance. Aerospace-ID Technologies White Paper Series. http://www.aero-id.org/research_reports/AEROID-CAM-010-TrackTrace.pdf (2006). Accessed 15 September 2008
[10] Folcke, G.: Traabilit des implants mdicaux en milieu hospitalier (in French). Revue de l'Electricit et de l'Electronique (REE), Socit Internationale des Electriciens (2008)
[11] Hardgrave, B., Patton, J.: RFID as electronic article surveillance EAS: feasibility assessment. Information Technology Research Institute Working Paper ITRI-WP117-0808. http://itrc.uark.edu/91.asp?code=completed&article=ITRI-WP117-0608 (2008). Accessed 25 August 2008
[12] Brusey, J., Floerkemeier, C., Harrison, M., Fletcher, M.: Reasoning about uncertainty in location identification with RFID. In: Proceedings of Workshop on Reasoning with Uncertainty in Robotics, IJCAI 03, Mexico (2003)
[13] Jeffery, S., Alonso, G., Franklin, M., Hong, W., Widom, J.: Declarative support for sensor data cleaning. In: Lecture Notes in Computer Science (LNCS) 3968/2006, 83-100. Springer, Heidelberg (2006)
[14] Jeffery, S., Garofalakis, M., Franklin, M.: Adaptive cleaning for RFID data streams. In: Proceedings of the 32nd International Conference on Very Large Data Bases (VLDB), 163-174. Korea (2006)
[15] Khoussainova, N., Balazinska, M., Suciu, D.: PEEX: Extracting Probabilistic Events from RFID Data. In: Proceedings of the 22th International Conference on Data Engineering (ICDE 08), 1480 - 1482, Mexico (2008)
[16] Kelepouris, T., McFarlane, D., Parlikad, A.: Developing a model for quantifying the quality and value of tracking information on supply chain decisions. In: Proceedings of the 12th International Conference on Information Quality (ICIQ-07), Boston (2007)
[17] EPCglobal Inc.: Pedigree Standard v. 1.0. EPCglobal Ratified Standard. http://www.epcglobalinc.org/standards/pedigree (2008). Accessed 14 September 2008
[18] Koh, R., Schuster, E., Chackrabarti, I., Bellman, A.: Securing the pharmaceutical supply chain. Auto-ID Labs White Paper. http://www.autoidlabs.org/uploads/media/MIT-AUTOID-WH021.pdf (2003). Accessed 14 September 2008
[19] Takaragi, K., Usami, M., Imura, R., Itsuki, R., Satoh, T.: An ultra small individual recognition security chip. IEEE Micro 21(6), 43-49 (2001)
[20] Lehtonen, M., Michahelles, F., Fleisch, E.: Probabilistic Approach for Location-Based Authentication. In: 1st International Workshop on Security for Spontaneous Interaction (IWSSI 07). 9th International Conference on Ubiquitous Computing, Austria (2007)
[21] EPCglobal Inc.: EPCIS 1.0.1 Specification. EPCglobal Ratified Standard. http://www.epcglobalinc.org (2007). Accessed 15 August 2008
[22] EPCglobal Inc.: EPCglobal Architecture Framework Version 1.2. www.epcglobalinc.org/ (2007). Accessed 25 August 2008
[23] John Jenkins Associates: Pharma traceability pilot – requirements analysis. Deliverable D6.2 of EU-BRIDGE Project. http://www.bridge-project.eu (2007).Accessed 15 September 2008