

Efficient Anonymous Category-level Joint Tag Estimation

Min Chen[†] Jia Liu[§] Shigang Chen[†] Qingjun Xiao[‡]

[†]Department of Computer & Information Science & Engineering
University of Florida, Gainesville, FL 32611, USA

[§]State Key Laboratory for Novel Software Technology
Nanjing University, P.R. China

[‡]Key Laboratory of Computer Network & Information Integration
Southeast University, P.R. China

Abstract—Radio-frequency identification (RFID) technologies have been widely used in many applications, including inventory management, supply chain, product tracking, transportation, logistics, etc. Tag estimation, which is to estimate the cardinality of a single tag set, is an important research topic. This paper expands the estimation research as follows: It performs joint estimation between two tag sets (which exist at different locations or at the same location but different times). More importantly the estimation is fine-grained in an effort to accommodate common practical scenarios, where each tag set consists of tags belonging to different categories. For any two given tag sets, we want to know the detailed information about the joint property of each category, instead of just the aggregate information of the whole sets. Furthermore, due to the open nature of RFID communications, it is often desirable that tag estimation can be performed in an anonymous way without revealing the tags' ID information. To support these requirements, we develop a new technique called mask bitmap that can encode a tag set without requiring the tags to report their IDs or category IDs. Any two mask bitmaps of different tag sets can be combined to perform category-level joint estimation. Through formal analysis, we determine how to set system parameters to meet a given accuracy requirement that can be arbitrarily set. Extensive simulation results confirm that the proposed solution can yield accurate category-level estimates in an efficient way, and preserve tags' anonymity as well.

I. INTRODUCTION

Radio Frequency Identification (RFID) technologies integrate simple communication, storage, and computation components into attachable tags that can communicate with RFID readers wirelessly over a distance [1], [2]. Due to this significant advantage over traditional bar code systems, RFID systems have been widely used in many applications [3]–[7]. Generally, an RFID system consists of three components: One or more readers, a large number of tags and a backend server. In RFID systems, tags with unique IDs are attached to objects, varying from products in a warehouse, merchandizes in a retail store, animals in a zoo, or medical equipments in a hospital. Each tag can not only identify the tagged object, but also indicate the category information through a subfield of tag ID called *category ID*.

One important RFID research topic is *tag estimation*, which is to estimate the *cardinality* of a tag set (i.e., the number of tags in the set) at a certain location [8]–[16]. Tag estimation can be used as a preprocessing step for optimizing the frame

size of frame-slotted ALOHA protocols in tag identification [10]. In addition, it can be applied to monitor the inventory level in a warehouse, the sales in a retail store, etc. Tag estimation is much more time-efficient than tag identification that needs to collect all tag IDs [10]. Moreover, since tag estimation does not need to identify any tags, the tags' anonymity can be preserved [9].

Although numerous approaches for tag estimation have been proposed, they have some limitations. First, most approaches only consider a single tag set [8]–[16]. Only limited prior work [17]–[20] studies the joint estimation of two or more tag sets. Moreover, almost all prior work, including [18]–[20], only estimates the aggregate information of the whole tag set(s), but ignores a common scenario where tagged objects may belong to different categories. Prior work [15], [16] investigates the category-level tag estimation, but they only consider a single tag set. This paper studies a new problem called category-joint tag estimation, which attempts to expand the research on tag estimation into a couple of new directions:

First, not only do we perform joint estimation between two tag sets, but more importantly the estimation is fine-grained at the category level in an effort to accommodate practical scenarios, where each tag set consists of tags belonging to different categories. Given two tag sets, we want to estimate the cardinality of the intersection set for each category. For example, consider a distribution network of warehouses, each carrying tagged products in various categories. For any two warehouses, if we have an automatic way to learn the cardinality of the intersection of their tag sets in each product category, we will have a means to track over time how each category of products flow through the network. For a single warehouse, knowing the cardinality of per-category intersection of tag sets at different times allows us to track how each category of products are moving in and out of the warehouse.

Second, we want to carry out category-level joint estimation anonymously without giving away the tags' private information, including tag IDs and category IDs [21], [22]. The widespread use of tags in traditional ways of deployment raises a serious privacy concern: The tags will report their IDs to any readers upon request, giving away the privacy of

the tag carriers. The RFID research community has recently devoted tremendous efforts to designing new mechanisms that keep the usefulness of tags while doing so anonymously [23]–[27], although they cannot be directly applied to the tag estimation problem. Leaking a category ID alone (without leaking the whole tag ID) will also be problematic. The category ID may indicate certain unique properties or private information about the tagged objects of a particular category. As an example, the category information of a medicine may identify its specific functions or targeted diseases, which can be used to infer what disease a patient may have if he/she buys this medicine. Therefore, we want to include the anonymity requirement in the design of a new protocol for category-level joint estimation. It is always preferred that a protocol can achieve the same functionality with comparable performance, but preserve the tags’ anonymity at the same time.

Category-level joint estimation is much more complicated than traditional tag estimation since we need to consider the joint properties of two tag sets, each further consisting of numerous different categories. Moreover, the anonymity requirement brings more challenges to the problem. Intuitively, each category of tags should be processed separately so that we can easily combine the information of a particular category from different tag sets. The difficulty is how to identify which category a tag belongs to without requiring it to report its category ID?

To our best knowledge, this is the first work that studies anonymous category-level joint estimation in RFID systems with new contributions summarized as follows:

First, we expand the traditional research on tag estimation — which only considers a single set or ignores the disparity of different tag categories — into new domains of category-level estimation and anonymity. The category-level joint estimation is capable of depicting the dynamics between two arbitrary tag sets at the category level.

Second, we propose a formal anonymous model to numerically evaluate the anonymity of different tag estimation protocols for the first time. In addition, we reveal the inherent tradeoff between estimation accuracy and anonymity of our proposed protocol.

Third, we develop a new technique called *mask bitmap* to achieve anonymous category-level joint estimation. Mask bitmaps can tactfully encode all tags of different categories without knowing their tag IDs or category IDs, while the information of each category can be retrieved later for joint estimation. We derive an estimator for the joint estimation of each category, formally analyze its mean and variance, and show that it can produce approximately unbiased results within an absolute error bound that can be set arbitrarily.

Finally, we perform extensive simulations to complement the theoretic analysis. The simulation results demonstrate that by following proper parameter settings our estimator can give very accurate estimates in an efficient way and preserve tags’ anonymity as well.

II. SYSTEM MODEL AND PROBLEM STATEMENT

A. System Model

Suppose all objects in an RFID system can be classified into m different categories, represented by a set M of category IDs $\{cid_1, cid_2, \dots, cid_m\}$. Each object is attached with a tag for identification purpose. Since each tag uniquely identifies one object, we will use object and tag interchangeably in the sequel. Given an arbitrary tag t , its tag ID id contains a subfield called category ID and is denoted as cid , indicating which specific category its associated object belongs to. Let a be the length of any tag ID (in number of bits) and $b (< a)$ be the length of any category ID.

RFID readers are installed to monitor tag sets located in their coverage areas. The readers can be connected to backend servers which provide supplemental storage, computation and communication resources. A reader communicates with tags using a frame-slotted ALOHA protocol. The reader initiates the communication by broadcasting a request that includes all necessary parameters. The tags are synchronized by the reader’s signal. In the following time frame, each tag sends its response in a randomly chosen slot. Based on the duration of slots, they can be classified into two types: One allows the transmission of a tag ID, whose duration is denoted by t_{id} ; the other is much shorter and only carries one bit information, whose duration is denoted by t_s . A t_s slot is called an *empty slot* if no tag replies, or a *busy slot* if one or multiple tags respond.

B. Anonymous Model

Low-cost RFID tags only have very limited computation, communication and storage resources. Hence, they cannot implement any classical cryptographic primitives, rendering the communications between a reader and a tag unprotected. Any adversary can plant unauthorized readers at chosen locations to eavesdrop on the transmissions between tags and readers, thereby capturing confidential or private information such as tag IDs and category IDs. We assume that the adversary has no prior knowledge of any tag IDs or category IDs in the system.

We use two numerical values to evaluate how much anonymity of a tag is preserved after executing a tag estimation protocol: (1) *ID anonymity*, which is the probability p_{id} that the adversary cannot infer a tag’s ID from the transmissions; (2) *Category anonymity*, which is the probability p_{cid} that the adversary cannot crack a tag’s category ID based on the transmissions.

C. Problem Statement

Consider two arbitrary tag sets N_p^* and N_q^* in a large distributed RFID system. They may refer to two sets at different locations, p and q , respectively, or two sets at the same location but different times, in which p and q refer to time. The wildcard superscript $*$ means that the notation covers tags of all categories. Let $n_p^* = |N_p^*|$ and $n_q^* = |N_q^*|$, $N_c^* = N_p^* \cap N_q^*$, and $n_c^* = |N_c^*|$, where the subscript c means “common tags” in the two sets N_p^* and N_q^* .

Given an arbitrary category $cid \in M$, $N_p^*(cid)$ is the subset of tags in N_p^* that belong to category cid , and

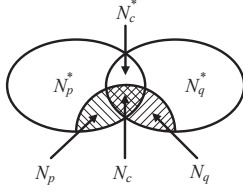


Fig. 1: The Venn diagram for sets N_p^* , N_q^* , N_c^* , N_p , N_q and N_c .

similarly $N_q(cid)$ is the subset of tags in N_q^* that belong to category cid . Let $n_p(cid) = |N_p(cid)|$, $n_q(cid) = |N_q(cid)|$, $N_u(cid) = N_p(cid) \cup N_q(cid)$, $n_u(cid) = |N_u(cid)|$, $N_c(cid) = N_p(cid) \cap N_q(cid)$, and $n_c(cid) = |N_c(cid)|$. Most of the time, our protocol description only needs to refer to one arbitrary category. Hence, we will abbreviate $N_p(cid)$ simply as N_p when the context does not raise ambiguity. Similarly, we will use N_q , N_u , N_c , n_p , n_q , n_u , and n_c without explicitly including (cid) in order to simplify these notations. The Venn diagram in Fig. 1 illustrates the relation between sets N_p^* , N_q^* , N_c^* , N_p , N_q and N_c .

The problem of anonymous category-level joint estimation is to estimate the value of n_c for each category under an accuracy requirement and an anonymity requirement. Once we have an estimate of the intersection cardinality n_c , it is trivial to estimate the union cardinality n_u and the difference cardinalities, $|N_p - N_q|$ and $|N_q - N_p|$, which are not presented in this paper due to space limitation. An intuitive interpretation for n_c is the number of tags in category cid that are transported from location p to location q or left behind from time p to time q if the two tag sets are recorded at the same location.

We use three performance metrics for evaluating anonymous category-level joint estimation, which are listed as follows.

Estimation accuracy: Our goal is to give an accurate estimate \hat{n}_c for n_c , such that

$$\text{Prob}\{|\hat{n}_c - n_c| \leq \eta\} \geq 1 - \theta, \quad (1)$$

where η is an absolute error bound and θ is a probability. For example, if $\eta = 50$ and $\theta = 10\%$, we require that the absolute estimation error $|\hat{n}_c - n_c|$ has a probability no less than 90% to be within $[0, 50]$. In other words, $[\hat{n}_c - \eta, \hat{n}_c + \eta]$ is a $(1 - \theta)$ confidence interval for n_c .

In contrast to our absolute error model, the traditional protocols [8]–[14], [16], [19], [20] for tag estimation employ a relative error model $\text{Prob}\{|\hat{n} - n| \leq \varepsilon n\} \geq 1 - \theta$, where n is the cardinality of a tag set, and ε is the relative error of the single-set estimation \hat{n} . This model has been adopted by the prior work on joint estimation of two tag sets [19], [20]. However, such adoption is not practically suitable in our opinion for the following reason: The single-set estimation protocols always assume a large tag set, which generally makes sense because a small set of tags does not need estimation — they can be directly counted. However, for two-set joint estimation, even though both sets are big, their intersection may be small or even empty (e.g., no object movement between two warehouses). According to the results in [20], the time complexity of the estimation procedure will approach to infinity when the intersection approaches to

empty, which is not acceptable because we cannot assume that the intersection of any two tag sets in a distributed RFID system is always large in practical applications. Hence, this paper advocates the absolute error model [18], which makes practical sense: For instance, consider a distribution network where each warehouse periodically encodes its tag set in an efficient, anonymous data structure (to be proposed later). We want to estimate the number of tagged products in each category that are moved between any two warehouses in each period, with an error of ± 50 tags at 95% confidence level. We will make the estimation by comparing the corresponding data structures.

Execution Time: Since RFID tags operate with low-speed communication channels, time efficiency is a key performance metric for all RFID protocols. Tag estimation should complete in a short time to avoid interference with other normal activities in an RFID system.

Anonymity: We use p_{id} and p_{cid} to measure the preserved anonymity of any tag after perform category-level joint estimation. More specifically, we want to maximize p_{id} and p_{cid} and make them as close to 1 as possible, such that it is practically infeasible for the adversary to infer the ID or category ID of any tag in N_p^* or N_q^* by eavesdropping on the execution of our estimation protocol.

III. RELATED WORK

There is no prior work on anonymous category-level joint estimation. Below we present some related work that can be applied to category-level joint estimation after slight modifications, and then point out their issues.

A. Tag Identification

The most straightforward approach for calculating n_c of an arbitrary category is to execute a tag identification protocol, e.g., Dynamic Framed Slotted ALOHA (DFSA) [28], [29] used by the EPC C1G2 standard. First, all tag IDs in N_p^* and N_q^* are collected by the readers. With the knowledge of N_p^* and N_q^* , we can easily find n_c for each category. One problem is that collecting all IDs is not efficient for large RFID systems, particularly when it has to be performed frequently. Due to transmission collisions, the lower bound of execution time to collect the IDs of n tags is $e \times n \times t_{id}$ [28], [29], where e is the natural constant. To make things worse, any IDs transmitted from the tags to the readers may be captured by the adversary, rendering $p_{id} = p_{cid} = 0$.

B. Prior Work on Joint Tag Estimation

Bitmap [30] is a compact data structure that is widely used for tag estimation [10], [31]. All bits in a bitmap are initialized to zeroes. Each tag is pseudo-randomly hashed to a bit in the bitmap and sets that bit to one. The ZDE protocol [19] leverages bitmap for joint estimation of two tag sets. The basic idea is to encode each tag set to an equal-size bitmap. Arbitrary two bitmaps can later be combined to estimate the cardinality of the union, differential, or intersection of the corresponding two tag sets. However, ZDE does not provide an explicit way of setting the bitmap size such that the estimation result can meet a given accuracy requirement. The

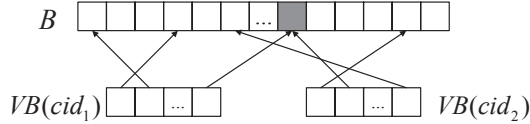


Fig. 2: Two virtual bitmaps $VB(cid_1)$ and $VB(cid_2)$ are built on top the mask bitmap B for categories cid_1 and cid_2 , respectively. The bit in grey is shared by both virtual bitmaps.

JREP protocol [18] is based on a similar idea, but improves the time efficiency of ZDE by enabling variable-size bitmaps. JREP first estimates the cardinality of a tag set, and then sets a proper size for the bitmap according to the estimated cardinality. Any two bitmaps can still be combined for joint tag estimation although they may have different sizes.

In [20], a generic Composite Counting Framework (CCF) is proposed to provide a cardinality estimate for any set expression with desired accuracy. CCF collects a synopsis for each tag set, which includes the d smallest hash values (of tag IDs) of that tag set. The synopses collected from different tag sets are leveraged to estimate the cardinality of a set expression.

To enable per-category estimation, we can employ ZDE, JREP, or CCF to perform joint estimation for every category. When performing estimation on a certain category, the reader can inform the tags belonging to that category to participate in the protocol execution, while other tags keep silent. However, the problem of ZDE and CCF is that their designs are based on a questionable relative error model, as we have explained previously. For example, CCF requires at least $\Theta(\frac{n_u}{\varepsilon^2 n_c} \ln \frac{1}{\theta})$ synopses to ensure $\text{Prob}\{|\hat{n}_c - n_c| \leq \varepsilon n_c\} \geq 1 - \theta$, which translates to excessively large execution time when $\frac{n_u}{n_c}$ is large, e.g., $n_c = 0$. JREP adopts an absolute error model, but it must first estimate the cardinality of each category to set a proper bitmap size. Although the execution time for estimating the cardinality of one category is small, the overall execution time can be very large when there are numerous categories. More importantly, all these approaches break the anonymity due to the transmissions of category IDs, resulting in $p_{cid} = 0$.

With the absolute error model, per-category estimation and the anonymity requirement, the prior work cannot solve this problem very well, which drives us to explore new ways for category-level joint estimation.

IV. ANONYMOUS CATEGORY-LEVEL JOINT ESTIMATION

A. Design Overview

Instead of performing joint estimation for each category one by one, which can be time-consuming, we want to enable category-level joint estimation in batch mode: All tags in one set, regardless of which categories they belong to, can be encoded to a bitmap simultaneously. Moreover, we want to avoid the transmissions of category IDs to protect the anonymity of tags. To achieve our objectives, we design a new data structure called *mask bitmap*, a variant of traditional bitmap. Our idea is to use a single large bitmap B to accommodate all categories. For each category, we build a *virtual bitmap* (VB) by randomly choosing some bits from

B , and any bit in B can be shared by multiple categories. Fig. 2 illustrates two virtual bitmaps $VB(cid_1)$ and $VB(cid_2)$ randomly chosen from the mask bitmap B for categories cid_1 and cid_2 , respectively, where the bit in grey is shared by both virtual bitmaps. A significant advantage of such bit-level sharing is that all categories use a common bitmap. Hence, each bit in the mask bitmap is shared by numerous tags in different categories, which helps conceal the tag ID and the category ID of a tag setting this bit.

Our anonymous category-level joint estimation protocol (CJEP) consists of two components: an encoding component for encoding a tag set to a mask bitmap, and an offline data analysis component to combine two arbitrary mask bitmaps, retrieve information of each category, and estimate n_c for each category (or any interested categories). Only the encoding component involves operations from the tags. In order to simplify the functions to be implemented on resource-constrained tags, our protocol follows an asymmetric design principle that pushes most complexity to the offline component while leaving the encoding component as simple as possible.

B. Encoding a Tag Set

We first describe the process of encoding a tag set covered by a single RFID reader. Consider the tag set N_p^* . We denote the mask bitmap for encoding N_p^* as B_p^* , consisting of f bits. Let $B_p^*[i]$ represent the i th bit in B_p^* , where $0 \leq i \leq f - 1$.

Virtual Bitmap: Each category cid is assigned a virtual bitmap $VB_p(cid)$ (abbreviated as VB_p) by pseudo-randomly taking l bits from B_p^* . This can be achieved by using l independent hash functions $H_k()$, each of which maps cid to the bit $B_p^*[H_k(cid)]$, where $0 \leq k \leq l - 1$ and the value of $H_k()$ is uniformly distributed over $[0, f)$. Denote the k th bit in VB_p as $VB_p[k]$. We have

$$VB_p[k] \equiv B_p^*[H_k(cid)]. \quad (2)$$

There are efficient ways to implement l hash functions on a tag. One approach is to employ a mater hash function $H^*()$ and a set R of l different random seeds as follows:

$$H_k(cid) = H^*(cid \oplus R[k]), \quad (3)$$

where \oplus is the XOR operator. The backend server generates the seeds and shares them with all readers, and each reader will include the seeds as part of the request message sent to the tags. Another approach relies on a pseudo-random number generator (which is required to be implemented on tags under the EPC C1G2 standard). Two random seeds, r_1 and r_2 , are needed. We use $cid \oplus r_1$ as the seed to the pseudo-random number generator. The value of $H_k(cid)$ is the $(kr_2 + 1)$ th output from the generator. That is, we use an output each time after dropping r_2 values.

Encoding: All bits in B_p^* are initialized to zeroes. The reader initiates the encoding process by broadcasting a request and the system parameters including the value of f and l random seeds. The request is followed by a time frame F , consisting of f slots. Consider an arbitrary tag with id as its ID and cid as its category ID. The purpose of the tag is to set a bit randomly chosen from the virtual bitmap of category cid to one. To do so, the tag will pseudo-randomly selects a

slot based on its own category ID and tag ID, which will be further explained shortly. The tag waits till the chosen slot to transmit a one-bit response. The reader keeps listening to the channel and sets $B_p^*[i] = 1$ if and only if the i th slot is busy.

There exists a one-to-one mapping between the i th slot in the time frame F and the i th bit in the mask bitmap B_p^* . Similarly, there is also a mapping from a virtual bitmap VB_p to the slots in F since all bits in VB_p are from B_p^* . The slots corresponding to the bits in VB_p form the category's *virtual frame*. Their indices in the whole frame F are $H_k(cid)$, where $0 \leq k < l$.

The tag will choose a slot from the virtual frame of category cid uniformly at random to transmit a one-bit response. To do that, it needs another hash function $h()$, which can also be implemented from the master hash function or the pseudo-random number generator as mentioned above, where the range of $h()$ is $[0, l)$. The tag chooses the $h(tid)$ th slot in the virtual time frame, corresponding to the $h(tid)$ th bit in VB_p . By (2), this bit is in fact the $H_{h(tid)}(cid)$ th bit in B_p^* , and in turn it corresponds to the $H_{h(tid)}(cid)$ th slot in the whole frame F . By the encoding design, the tag effectively sets the following bit to one.

$$VB_p[h(tid)] \equiv B_p^*[H_{h(tid)}(cid)] = 1. \quad (4)$$

We stress that the tag never constructs the virtual bitmap VB_p explicitly. Its operation is actually very simple: All it does is to compute two hashes for the value of $H_{h(tid)}(cid)$ and then waits for that slot to transmit a signal, which sets a randomly chosen bit in the virtual bitmap of category cid to one.

Multiple readers: In case that multiple readers are needed to cover all tags in the system, we assume that a reader schedule is established based on signal measurement in order to avoid reader-to-reader collisions. Only non-conflicting readers, each covering a non-overlapping area, are scheduled to be active at the same time. After every reader reports its mask bitmap, the backend server performs a bitwise OR operation over all mask bitmaps to obtain a combined mask bitmap that encodes the whole tag set. The prescribed system parameters, including the time frame size f and virtual bitmap size l , are used by all readers across the system. We will introduce an algorithm for determining the system parameters shortly.

C. Offline Information Retrieval

After N_p^* and N_q^* are encoded to mask bitmaps B_p^* and B_q^* , respectively, the bitmaps are offloaded to the backend server for permanent storage. For a query to estimate the cardinality n_c of their intersection on an arbitrary category cid , the backend server first retrieves the two virtual bitmaps of category cid from B_p^* and B_q^* as follows:

$$\begin{aligned} VB_p[k] &\equiv B_p^*[H_k(cid)] \\ VB_q[k] &\equiv B_q^*[H_k(cid)], \end{aligned} \quad (5)$$

where $0 \leq k \leq l-1$. The server then performs a bitwise OR operation on the two virtual bitmaps, resulting in another l -bit bitmap, denoted by VB_u ,

$$VB_u[k] = VB_p[k] \vee VB_q[k], \quad 0 \leq k \leq l-1, \quad (6)$$

where the subscript u means that VB_u encodes the ‘‘union’’ of $N_p(cid)$ and $N_q(cid)$. Fig. 3 shows the process of retrieving VB_p , VB_q , and generating VB_u from B_p^* and B_q^* . By the

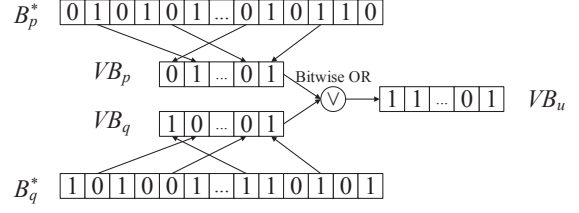


Fig. 3: The process of retrieving VB_p , VB_q , and generating VB_u from B_p^* and B_q^* .

encoding design, all tags in N_p are encoded in VB_p in a probabilistic way; recall that N_p is abbreviation of $N_p(cid)$. Similarly, all tags in N_q are encoded in VB_q . Hence, all tags in $N_p \cup N_q$ are encoded in VB_u .

However, as we have discussed previously and illustrated in Fig. 2, different virtual bitmaps share bits, which means that a bit of ‘1’ in VB_p (or VB_q , VB_u) may not be set by tags in category cid , but instead set by tags in another category. Hence, while bit sharing helps achieve anonymity and efficiency, it also introduces inter-category noise. To deal with the noise issue and accurately estimate the value of n_c , we resort to probabilistic analysis in the next subsection.

In addition, the backend server combines B_p^* and B_q^* by performing a bitwise OR operation to obtain another useful bitmap, denoted by B_u^* . Hence,

$$B_u^*[k] = B_p^*[k] \vee B_q^*[k], \quad 0 \leq k \leq f-1. \quad (7)$$

We know that B_p^* , B_q^* and B_u^* encode the tag sets, N_p^* , N_q^* and $N_p^* \cup N_q^*$, respectively.

D. Estimator for n_c

In this Section, we derive an estimator \hat{n}_c for estimating n_c using B_p^* , B_q^* , B_u^* , VB_p , VB_q and VB_u . We first present and prove the following theorem.

Theorem 1. *An arbitrary tag t has a probability $\frac{1}{f}$ to be mapped to a given bit z in a f -bit mask bitmap.*

Proof: The process for a tag to randomly choose a slot in a time frame can be cast into bins and balls problem [32]. We denote the l -bit virtual bitmap for t 's category as VB . Let random variable X represent the number of physical bits occupied by VB in the mask bitmap. We have

$$Prob(X = x) = \frac{\binom{f}{x} \times x! \times S(l, x)}{f^l},$$

where $S(l, x) = \frac{1}{x!} \sum_{i=0}^x (-1)^i \binom{x}{i} (x-i)^l$ is the Stirling number of the second kind [33]. $S(l, x)$ gives the number of ways to partition a set of l balls into x non-empty bins. In addition, the probability that z is one of the x bits is $1 - \frac{\binom{f-1}{x}}{\binom{f}{x}}$. The tag ID has a probability $\frac{1}{x}$ to be mapped to z among those x bits. Therefore, the probability for t to be mapped to z is

$$\begin{aligned} p_z &= \sum_{x=1}^l \frac{\binom{f}{x} \times x! \times S(l, x)}{f^l} \times \left(1 - \frac{\binom{f-1}{x}}{\binom{f}{x}}\right) \times \frac{1}{x} \\ &= \frac{1}{f^{l+1}} \times \sum_{x=1}^l \binom{f}{x} \times x! \times S(l, x) = \frac{1}{f}, \end{aligned} \quad (8)$$

where we have used $\sum_{x=1}^l \binom{f}{x} \times x! \times S(l, x) = f^l$ [33]. ■

Now let us continue to derive an estimator \hat{n}_c . Let X_j be the event that the j th bit in B_p^* is 0 ($0 \leq j \leq l-1$), and 1_{X_j} be the corresponding indicator random variable, namely,

$$1_{X_j} = \begin{cases} 1, & \text{if } B_p^*[j] = 0, \\ 0, & \text{if } B_p^*[j] = 1. \end{cases}$$

Therefore, we have $Prob(X_j) = (1 - \frac{1}{f})^{n_p^*}$, and $E(1_{X_j}) = 1 \times Prob(X_j) + 0 \times (1 - Prob(X_j)) = (1 - \frac{1}{f})^{n_p^*}$. Let U_p be a random variable of the fraction of bits in B_p^* that remain zeros after encoding all tags in N_p^* . We have $U_p = \frac{1}{f} \sum_{j=0}^{f-1} 1_{X_j}$. Hence,

$$E(U_p) = \frac{1}{f} \sum_{j=0}^{f-1} E(1_{X_j}) = (1 - \frac{1}{f})^{n_p^*}. \quad (9)$$

Similarly, let Y_j be the event that the j th bit in B_q^* is 0, and U_q be a random variable for the fraction of bits in B_q^* that remain zeros. We have

$$E(U_q) = (1 - \frac{1}{f})^{n_q^*}. \quad (10)$$

Let Z_j be the event that the j th bit in B_u^* is 0. Since this bit is OR of the j th bit in B_p^* and B_q^* , $Prob(Z_j) = Prob(X_j \wedge Y_j) = Prob(X_j|Y_j) \times Prob(Y_j) = (1 - \frac{1}{f})^{n_p^* + n_q^* - n_c^*}$. Let random variable U_u denote the fraction of zeroes in B_u^* . We have

$$E(U_u) = (1 - \frac{1}{f})^{n_p^* + n_q^* - n_c^*}. \quad (11)$$

Combining (9), (10) and (11), we know $E(U_u) = E(U_p)E(U_q)(1 - \frac{1}{f})^{-n_c^*}$. Therefore,

$$(1 - \frac{1}{f})^{-n_c^*} = \frac{E(U_u)}{E(U_p)E(U_q)}. \quad (12)$$

Now let us move forward to investigate the properties of VB_p , VB_q and VB_u . Let C_j ($0 \leq j \leq l-1$) be the event that the j th bit in VB_p is 0, and 1_{C_j} be the corresponding indicator random variable. For event C_j to happen, neither a tag in N_p nor a tag in $N_p^* - N_p$ shall be mapped to $VB_p[j]$. The probability is $Prob(C_j) = (1 - \frac{1}{l})^{n_p} (1 - \frac{1}{f})^{n_p^* - n_p}$, and therefore $E(1_{C_j}) = (1 - \frac{1}{l})^{n_p} (1 - \frac{1}{f})^{n_p^* - n_p}$. Let V_p represent the fraction of 0s in VB_p . We have

$$E(V_p) = \frac{1}{l} \sum_{j=0}^{l-1} E(1_{C_j}) = (1 - \frac{1}{l})^{n_p} (1 - \frac{1}{f})^{n_p^* - n_p}. \quad (13)$$

Applying (9) to (13), we have

$$E(V_p) = (1 - \frac{1}{l})^{n_p} (1 - \frac{1}{f})^{-n_p} E(U_p). \quad (14)$$

Substituting $E(V_p)$, $E(U_p)$ with V_p , U_p , respectively, and taking the logarithm of the both sides, we derive an estimator for n_p as follows:

$$\hat{n}_p = \frac{\ln V_p - \ln U_p}{\ln(1 - \frac{1}{l}) - \ln(1 - \frac{1}{f})}. \quad (15)$$

Let D_j be the event that the j th bit in VB_q is 0, and V_q represent the fraction of 0s in VB_q . Similarly, we have

$$E(V_q) = (1 - \frac{1}{l})^{n_q} (1 - \frac{1}{f})^{n_q^* - n_q}, \quad (16)$$

$$\hat{n}_q = \frac{\ln V_q - \ln U_q}{\ln(1 - \frac{1}{l}) - \ln(1 - \frac{1}{f})}. \quad (17)$$

Consider an arbitrary bit z in VB_u . Only if the following two conditions are satisfied will its value remains zero.

- 1) z is not chosen by any tag in $N_p \cup N_q$.
- 2) z is not chosen by any tag in $(N_p^* \cup N_q^*) - (N_p \cup N_q)$.

For the first condition, each tag in $N_p \cup N_q$ has a probability $\frac{1}{f}$ to select z and set it to 1. Hence, the probability q_1 to satisfy condition one can be calculated by

$$q_1 = (1 - \frac{1}{f})^{n_p + n_q - n_c}. \quad (18)$$

For the second condition, each tag in $(N_p^* \cup N_q^*) - (N_p \cup N_q)$ has a probability $\frac{1}{f}$ to choose z . Hence,

$$q_2 = (1 - \frac{1}{f})^{n_p^* + n_q^* - n_c^* - (n_p + n_q - n_c)}. \quad (19)$$

Let E_j be the event that the j th bit in VB_u is 0, 1_{E_j} be the corresponding indicator random variable, and V_u be the fraction of 0s in VB_u . Combining (18) and (19), we have

$$Prob(E_j) = (1 - \frac{1}{l})^{n_p + n_q - n_c} (1 - \frac{1}{f})^{n_p^* + n_q^* - n_c^* - (n_p + n_q - n_c)},$$

$$E(V_u) = (1 - \frac{1}{l})^{n_p + n_q - n_c} (1 - \frac{1}{f})^{n_p^* + n_q^* - n_c^* - (n_p + n_q - n_c)}. \quad (20)$$

We know $n_u = n_p + n_q - n_c$. Applying (11) to (20),

$$E(V_u) = (1 - \frac{1}{l})^{n_u} (1 - \frac{1}{f})^{-n_u} E(U_u).$$

$$\hat{n}_u = \frac{\ln V_u - \ln U_u}{\ln(1 - \frac{1}{l}) - \ln(1 - \frac{1}{f})}. \quad (21)$$

Applying (12), (13) and (16) to (20), we have

$$E(V_u) = E(V_p)E(V_q) \left(\frac{1 - \frac{1}{f}}{1 - \frac{1}{l}} \right)^{n_c} \frac{E(U_u)}{E(U_p)E(U_q)}.$$

Hence, we can calculate

$$n_c = \frac{\ln \frac{E(V_u)}{E(V_p)E(V_q)} - \ln \frac{E(U_u)}{E(U_p)E(U_q)}}{\ln(1 - \frac{1}{f}) - \ln(1 - \frac{1}{l})}. \quad (22)$$

Replacing $E(V_u)$, $E(V_p)$, $E(V_q)$, $E(U_u)$, $E(U_p)$ and $E(U_q)$ with observed values V_u , V_p , V_q , U_u , U_p and U_q , respectively, we obtain an estimator for n_c as follows

$$\hat{n}_c = \frac{(\ln V_p - \ln U_p) + (\ln V_q - \ln U_q) - (\ln V_u - \ln U_u)}{\ln(1 - \frac{1}{l}) - \ln(1 - \frac{1}{f})}. \quad (23)$$

Note that f and l should be set properly such that U_p , U_q , U_u , V_p , V_q , and V_u are non-zero.

E. Analysis of \hat{n}_c

We have derived $E(\hat{n}_c)$ and $Var(\hat{n}_c)$ but cannot present the derivation process here due to space limitation.

The expected value of \hat{n}_c is

$$E(\hat{n}_c) \approx n_c + \frac{e^{v_p + w_p} + e^{v_q + w_q} - e^{v_u + w_u} - v_p - v_q + v_u - 1}{2}. \quad (24)$$

Note that when the values of v_p , w_p , v_q , w_q , v_u and w_u are small, $Bias(\hat{n}_c) = E(\hat{n}_c) - n_c \approx 0$, which means \hat{n}_c is an asymptotically unbiased estimator of n_c .

Since the close form of $Var(\hat{n}_c)$ is extremely complicated, we also obtain an upper bound for $Var(\hat{n}_c)$ that takes a much simpler form as follows:

$$Var(\hat{n}_c) < \frac{e^{v_u + w_u} - v_u - 1}{c^2 l} + \frac{2l(1 - e^{-v_u})}{c^2 f}, \quad (25)$$

where $c = \ln(1 - \frac{1}{l}) - \ln(1 - \frac{1}{f})$.

Protocol	Identification	ZDE/JREP	CCF	CJEP
p_{cid}	0	0	0	$1 - \frac{f}{l2^b}$
p_{id}	0	$1 - \frac{f}{2^{a-b}}$	$1 - \frac{D}{2^{a-b}}$	$1 - \frac{f}{2^a}$

TABLE I: Preserved anonymity of a tag after executing different protocols for category-level joint estimation, where a is the length of tag IDs and b is the length of category IDs.

F. Analysis of Anonymity

In this section, we analyze the preserved anonymity of a tag after executing CJEP under our proposed anonymous model. We assume that the adversary has unlimited computing and storage resources such that given an arbitrary slot it knows which tags will be mapped to this slot (which requires $O(2^a)$ preprocessing overhead). This assumption overrates the adversary's capability, so the following two theorems provide a lower bound of the anonymity of CJEP.

Theorem 2. *Given an arbitrary tag and its response in a time frame including f slots, $p_{cid} = 1 - \frac{f}{l2^b}$.*

Proof: With a b -bit category ID, there are 2^b different category IDs. For each category, a l -bit virtual bitmap is allocated, which corresponds to l different slots. Hence, the mean number of categories mapped to each slot is $\frac{l2^b}{f}$, and there is no other clue for the adversary to distinguish any two categories mapped to the same slot. As a result, the probability for the adversary to guess the category ID of a tag based on its response is $\frac{1}{\frac{l2^b}{f}} = \frac{f}{l2^b}$. Hence, $p_{cid} = 1 - \frac{f}{l2^b}$. ■

Theorem 3. *Given an arbitrary tag and its response in a time frame including f slots, $p_{id} = 1 - \frac{f}{2^a}$.*

Proof: According to Theorem 2, the adversary has a probability $\frac{f}{l2^b}$ to correctly guess its category ID based on its chosen slot. With a a -bit tag ID and a b -bit category ID, each category can have as many as 2^{a-b} tags, which are evenly distributed to l slots (l -bit virtual bitmap). Therefore, each slot averagely has $\frac{2^{a-b}}{l}$ tags belonging to that category. Therefore, the probability to infer the full ID correctly is $\frac{f}{l2^b} \times \frac{l}{2^{a-b}} = \frac{f}{2^a}$. Hence, $p_{id} = 1 - \frac{f}{2^a}$, which is extremely small when reasonably long tag IDs, e.g., 120 bits, are used. ■

Table I compares the preserved anonymity of a tag after executing different protocols for category-level joint estimation. Only CJEP can preserve category anonymity. For CCF [20], D is the size of hash space, and it should be set to $O(n^2)$, where n is the cardinality of a tag set. Generally, we have $f < D$. Therefore, CJEP has the best ID anonymity when the same frame size is used by ZDE/JREP.

G. Parameter Setting

In this section, we propose an algorithm to set appropriate parameters f and l under the accuracy requirement given in (1). We use $g(f, l)$ to represent the upper bound of $Var(\hat{n}_c)$

Algorithm 1 Procedure of determining optimal f and l .

Input: $n_p^*, n_q^*, n_u^*, n_p, n_q, n_u, \eta, \theta$

1: $f = 1$ {Initializes f }

2: **repeat**

3: $f = f + s$ { s is the step size for increasing f }

4: $w_p = \frac{n_p^*}{f}, w_q = \frac{n_q^*}{f}, w_u = \frac{n_u^*}{f}$

5: binary search for l^* that satisfies $\frac{\partial g}{\partial l} = 0$

6: $v_p = \frac{n_p}{l}, v_q = \frac{n_q}{l}, v_u = \frac{n_u}{l}$

7: **until** $Var(\hat{n}_c) \leq \eta^2 / Z_{\frac{\theta}{2}}^2$

Output: f, l

given in (25):

$$g(f, l) = \frac{e^{v_u+w_u} - v_u - 1}{c^2 l} + \frac{2l(1 - e^{-v_u})}{c^2 f} \quad (26)$$

$$\approx l(e^{v_u+w_u} - v_u - 1) + \frac{2l^3(1 - e^{-v_u})}{f},$$

since $c^2 = (\ln(1 - \frac{1}{l}) - \ln(1 - \frac{1}{f}))^2 \approx (-\frac{1}{l} + \frac{1}{f})^2 \approx \frac{1}{l^2}$ when $l \ll f$. Taking the partial derivative of $g(f, l)$ with respect to f , we have

$$\frac{\partial g}{\partial f} = -\frac{lw_u e^{v_u+w_u}}{f} - \frac{2l^3(1 - e^{-v_u})}{f^2} < 0.$$

Hence, $g(f, l)$ decreases with the increase of f . For the setting of l , we have the following theorem.

Theorem 4. *Given a value of f , there exists one and only one optimal l , denoted by l^* , that minimizes $g(f, l)$.*

Proof: Taking the partial derivative of $g(f, l)$ with respect to l , we have

$$\frac{\partial g}{\partial l} = e^{v_u+w_u}(1 - v_u) - 1 + \frac{2l^2}{f}(3(1 - e^{-v_u}) - v_u e^{-v_u}).$$

In addition, we can calculate

$$\frac{\partial^2 g}{\partial l^2} = \frac{v_u^2 e^{v_u+w_u}}{l} + \frac{2l}{f}(6 - 6e^{-v_u} - 4v_u e^{-v_u} - v_u^2 e^{-v_u}).$$

It is easy to prove that $(6 - 6e^{-v_u} - 4v_u e^{-v_u} - v_u^2 e^{-v_u})$ is monotonically increasing with respect to v_u . Hence $(6 - 6e^{-v_u} - 4v_u e^{-v_u} - v_u^2 e^{-v_u}) \geq 0$, with the minimum reaching at $v_u = 0$. Therefore, $\frac{\partial^2 g}{\partial l^2} \geq 0$, which implies that $\frac{\partial g}{\partial l}$ is monotonically increasing with respect to l . Moreover, we have $\frac{\partial g}{\partial l}|_{l=0} \leq 0$ and $\frac{\partial g}{\partial l}|_{l \rightarrow \infty} > 0$. Therefore, there is only one value l^* that satisfies $\frac{\partial g}{\partial l} = 0$, which minimizes $g(f, l)$. In other words, there is an optimal value of l that minimizes the upper bound of $Var(\hat{n}_c)$ given an arbitrary setting of f . ■

For a normal distribution with $E(\hat{n}_c) \approx n_c$, the requirement in (1) can be translated to

$$Z_{\frac{\theta}{2}} \sqrt{Var(\hat{n}_c)} \leq \eta,$$

where $Z_{\frac{\theta}{2}}$ is the $1 - \frac{\theta}{2}$ percentile for the standard Normal distribution. Therefore,

$$Var(\hat{n}_c) \leq \eta^2 / Z_{\frac{\theta}{2}}^2. \quad (27)$$

We need to set f large enough such that $Var(\hat{n}_c)$ is no larger than $\eta^2 / Z_{\frac{\theta}{2}}^2$. For a given value of f , we calculate the optimal l that minimizes $g(f, l)$. With this (f, l) pair, we check whether the condition in (27) satisfies. If not, we further

increase the value of f to reduce $Var(\hat{n}_c)$. Algorithm 1 shows the procedure of determining optimal f and l . In practice, the values of n_p^* , n_q^* , n_u^* , n_p , n_q , and n_u are unknown. We will shortly show how to set those parameters in Section V.

V. SIMULATION RESULTS

In this section, we evaluate the performance of CJEP through simulations. As we have explained in Section III, there is no prior work on anonymous category-level joint estimation. Therefore, we use a tag identification protocol as a benchmark for comparison. In addition, we apply state-of-the-art protocols on joint tag estimation, which are JREP [18] and CCF [20], to per-category joint estimation after some slight modifications.

A. Parameter Settings

In our simulations, the communication parameters are set following the specification of EPC C1G2 standard [1]. Any two consecutive communications between the reader and tags are separated by a time interval of $302 \mu s$. The transmission rate between the reader and tags is in the range of $26.7 kbps$ to $128 kbps$. We set the transmission rate to $26.7 kbps$ (similar simulation results can be observed under other parameter configurations). That is, it takes the reader or tags $37.45 \mu s$ to transmit one bit. Therefore, we have $t_s = 37.45 + 302 = 339.45 \mu s$, and $tid = 37.45 \times 96 + 302 = 3897.2 \mu s$.

We set the number m of categories to 1000 to ensure there are enough categories for evaluating CJEP in a probabilistic way. We let the value n_c of each category follow a uniform distribution over $[0, 500]$, so that every value of n_c have enough samples. In fact, the adoption of absolute error model guarantees that CJEP works regardless of the distribution of n_c . We let the numbers n_p and n_q of tags in each category follow several common distributions as follows:

Dist. 1: n_p and n_q independently follow a uniform distribution $Unif(300, 700)$.

Dist. 2: n_p and n_q independently follow a normal distribution $Norm(500, 250^2)$.

Dist. 3: n_p and n_q of independently follow a zipf distribution [34] over $[400, 1000]$ at steps of 10 (61 different values in total). With the value of the exponent characterizing the distribution set to 1.0, the frequency of the rank- j value is $f(j) = \frac{1/j}{\sum_{i=1}^{61} (1/i)}$.

The absolute error bound η varies from 20 to 100, at steps of 10. In addition, we set θ to 0.1 and 0.05, and $Z_{\frac{\theta}{2}}$ is therefore 1.645 and 1.960, respectively. The values of f and l are obtained from Algorithm 1, where the step size s is set to 1000. We set $n_p^* = n_q^* = 500000$, which gives an estimated upper bound of the cardinality of a tag set¹ (e.g., the number of goods can be stored in a warehouse). In addition, we approximately set $n_p = \frac{n_p^*}{m}$, $n_q = \frac{n_q^*}{m}$, $n_u = n_p^* + n_q^*$, and $n_u = n_p + n_q$.

For CCF, the size D of hash space is set to 100000 to avoid hash collisions. Hence, the length of each hash value is

¹Note that the purpose for the large setting of n_p^* and n_q^* here is to ensure there are enough categories and each category contains sufficient number of tags for evaluating CJEP in a probabilistic way.

$\lceil \log_2 100000 \rceil = 17$ bits. Instead of setting the number d of synopses to $\Theta(\frac{n_u}{\varepsilon^2 n_c} \ln \frac{1}{\theta})$, which otherwise can be excessively large when n_c is small (n_c is not known in advance), we fix d to 400. The setting of JREP exactly follows that in [18].

B. Execution time

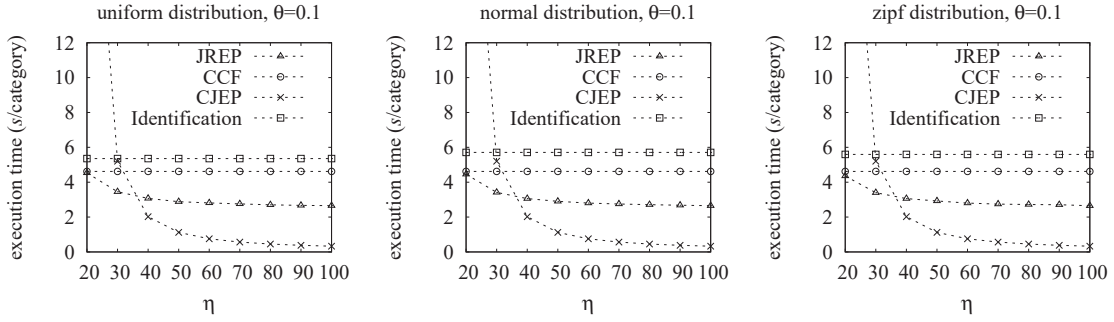
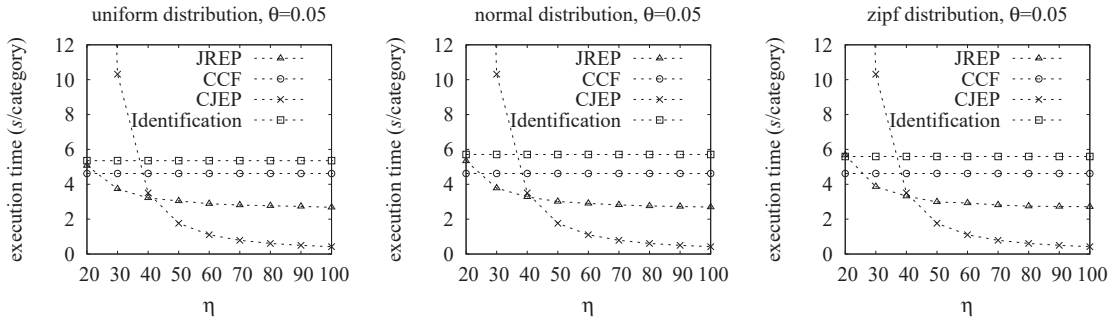
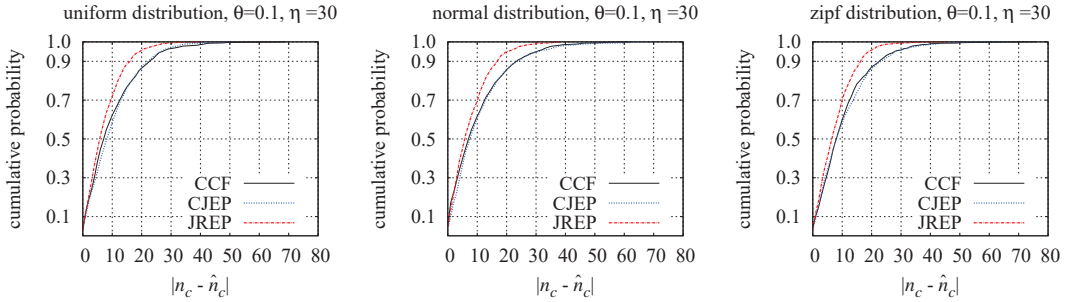
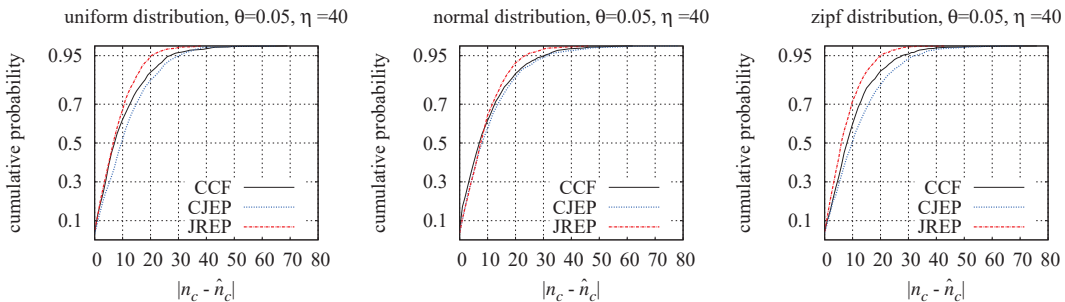
We first compare the execution time of different protocols for category-level joint estimation. Fig. 4 and Fig. 5 show the results when $\theta = 0.1$ and $\theta = 0.05$, respectively. In each plot, the x axis is the absolute error bound η , and the y axis is the average execution time of each category. Similar results can be observed even though n_p and n_q follow different distributions. This is because the total number of tags in each tag set is close and the number of categories is the same in different settings. Both JREP and CJEP take less execution time with the increase of η , meaning a larger estimation error is allowed. The difference is that the execution time of CJEP decreases more dramatically than that of JREP with the increase of η . Although CJEP takes a longer execution time than the other three protocols when η is very small, it is much more efficient than others when η is moderately large. For example, when $\eta = 50$ and $\theta = 0.1$, CJEP only takes 38.7%, 24.3%, 19.7% of the execution time of JREP, CCF and an identification protocol, respectively. In addition, a smaller θ requires that the absolute estimation error $|n_c - \hat{n}_c|$ has a higher probability to be bounded by η (i.e., a more stringent requirement on estimation accuracy). Therefore, the execution time of CJEP increases when θ decreases to 0.05 from 0.1.

C. Estimation accuracy

Next we evaluate the estimation accuracy of different protocols. For the purpose of fair comparison, we choose a setting of η and θ that makes the protocols have close execution times. More specifically, we set $\eta = 40$ when $\theta = 0.1$, and $\eta = 30$ when $\theta = 0.05$. Since an identification protocol has no estimation error for n_c , we do not include the identification protocol in the performance comparison. For the (n_c, \hat{n}_c) pair of each category in the estimation results, we calculate the absolute estimation error $|n_c - \hat{n}_c|$. Fig. 6 and Fig. 7 show the cumulative distribution function (CDF) of the absolute estimation error. In Fig. 6, the 90 percentile of $|n_c - \hat{n}_c|$ of JREP, CCF, and CJEP are respectively 16, 22, and 23, which are all within the error bound $\eta = 30$. In Fig. 7, the 95 percentile of $|n_c - \hat{n}_c|$ of JREP, CCF, and CJEP are respectively 20, 26, and 30, which are also bounded by $\eta = 40$. Overall, the estimation accuracy of JREP is a little better than that of CCF and CJEP, and CCF and CJEP yield estimates with comparable accuracy. The bit-level sharing in mask bitmaps of CJEP helps preserve tags' anonymity, but also introduces inter-category noise that brings some negative effect on the estimation accuracy.

D. Anonymity

Recall that an identification protocol does not preserve any ID anonymity or category anonymity. From Table I, we know $p_{id} \approx 1$ for CCF, JREP and CJEP when the ID length a is large enough, e.g., 96 bits. Since CCF and JREP cannot


 Fig. 4: Execution time comparison of different protocols when $\theta = 0.1$.

 Fig. 5: Execution time comparison of different protocols when $\theta = 0.05$.

 Fig. 6: Cumulative probability of the the absolute estimation error $|n_c - \hat{n}_c|$ when $\eta = 30$ and $\theta = 0.1$.

 Fig. 7: Cumulative probability of the the absolute estimation error $|n_c - \hat{n}_c|$ when $\eta = 40$ and $\theta = 0.05$.

$\theta \backslash \eta$	20	30	40	50	60	70	80	90	100
0.1	96.69%	99.01%	99.51%	99.69%	99.78%	99.82%	99.85%	99.88%	99.89%
0.05	94.03%	98.38%	99.26%	99.56%	99.69%	99.77%	99.81%	99.84%	99.86%

 TABLE II: p_{cid} of CJEP when $b = 20$.

$\eta \backslash \theta$	20	30	40	50	60	70	80	90	100
0.1	99.90%	99.97%	99.98%	99.99%	99.99%	99.99%	100.00%	100.00%	100.00%
0.05	99.81%	99.95%	99.98%	99.99%	99.99%	99.99%	99.99%	100.00%	100.00%

TABLE III: p_{cid} of CJEP when $b = 25$.

preserve category anonymity, we focus on investigating the p_{cid} of CJEP.

Table II and Table III show the values of p_{cid} when the length b of category ID is set to 20 bits and 25 bits, respectively. Under the same setting, a larger value of b , which makes the category ID more difficult to guess, leads to a larger value of p_{cid} . When $b = 25$, p_{cid} approaches to 1 in most cases. In addition, we observe that with fixed values of b and θ , p_{cid} increases with the increase of η ; with fixed values of b and η , p_{cid} increases with the increase of θ . Generally speaking, a smaller value of η or θ means a higher requirement of estimation accuracy, which requires a longer execution time (a larger f). According to Theorem 2 and Theorem 3, the increase of f causes p_{cid} and p_{id} to decrease. Therefore, there exists a tradeoff between estimation accuracy and preserved anonymity.

VI. CONCLUSION

This paper studies a new problem of anonymous category-level joint estimation in RFID systems: Given a particular category, we want to estimate its cardinality in the intersection of two arbitrary tag sets anonymously. We propose a protocol CJEP based on a novel data structure called mask bitmap. We derive an estimator, analyze its mean and variance, and provide an algorithm for system parameter setting. We also point out the inherent tradeoff between the estimation accuracy and anonymity using CJEP. We perform extensive simulations to evaluate the performance of our protocol.

VII. ACKNOWLEDGMENTS

This work was supported in part by the National Science Foundation under grant CNS-1409797 .

REFERENCES

- [1] *EPC Radio-Frequency Identity Protocols Class-1 Gen-2 UHF RFID Protocol for Communications at 860MHz-960MHz*, EPCglobal, Available at <http://www.epcglobalinc.org/uhfclg2>.
- [2] J. Wang, H. Hassanieh, D. Katabi, and P. Indyk, "Efficient and Reliable Low-power Backscatter Networks," *Proc. of ACM SIGCOMM*, 2012.
- [3] M. Chen, W. Luo, Z. Mo, S. Chen, and Y. Fang, "An efficient tag search protocol in large-scale rfid systems with noisy channel," *IEEE/ACM Transactions on Networking*, vol. 24, no. 2, pp. 703–716, April 2016.
- [4] J. Liu, M. Chen, B. Xiao, F. Zhu, S. Chen, and L. Chen, "Efficient rfid grouping protocols," *IEEE/ACM Transactions on Networking*, vol. PP, no. 99, pp. 1–1, 2016.
- [5] S. Qi, Y. Zheng, M. Li, Y. Liu, and J. Qiu, "Scalable Data Access Control in RFID-Enabled Supply Chain," *Proc. of IEEE ICNP*, pp. 71–82, O 2014.
- [6] F. Li, "Impact of RFID Technology on Supply Chain: A Simulation Approach," *Proc. of IEEE MASS*, pp. 1–4, September 2009.
- [7] C.H. LEE and C.W. Chung, "Efficient Storage Scheme and Query Processing for Supply Chain Management using RFID," *Proc. of ACM SIGMOD*, pp. 291–302, October 2008.
- [8] B. Chen, Z. Zhou, and H. Yu, "Understanding RFID Counting Protocols," *Proc. of ACM MOBICOM*, pp. 291–302, 2013.
- [9] H. Han, B. Sheng, C. Tan, Q. Li, W. Mao, and S. Lu, "Counting RFID Tags Efficiently and Anonymously," *Proc. of IEEE INFOCOM*, 2010.
- [10] M. Kodialam and T. Nandagopal, "Fast and Reliable Estimation Schemes in RFID Systems," *Proc. of ACM MOBICOM*, 2006.
- [11] T. Li, S. Wu, S. Chen, and M. Yang, "Energy Efficient Algorithms for the RFID Estimation Problem," *Proc. of IEEE INFOCOM*, March 2010.
- [12] M. Shahzad and A. Liu, "Every Bit Counts: Fast and Scalable RFID Estimation," *Proc. of ACM MOBICOM*, pp. 365–376, August 2012.
- [13] Y. Zheng and M. Li, "ZOE: Fast Cardinality Estimation for Large-scale RFID Systems," *Proc. of IEEE INFOCOM*, pp. 908–916, April 2013.
- [14] Y. Zheng, M. Li, and C. Qian, "PET: Probabilistic Estimating Tree for Large-Scale RFID Estimation," *Proc. of IEEE ICDCS*, 2011.
- [15] W. Luo, Y. Qiao, and S. Chen, "An Efficient Protocol for RFID Multigroup Threshold-based Classification," *Proc. of IEEE INFOCOM*, pp. 890–898, April 2013.
- [16] L. Xie, H. Han, Q. Li, J. Wu, and S. Lu, "Efficiently Collecting Histograms Over RFID Tags," *Proc. of IEEE INFOCOM*, 2014.
- [17] Qingjun Xiao, Shigang Chen, and Min Chen, "Joint property estimation for multiple rfid tag sets using snapshots of variable lengths," *Proc. of ACM Mobihoc*, pp. 151–160, 2016.
- [18] Qingjun Xiao, Min Chen, Shigang Chen, and Yian Zhou, "Temporally or spatially dispersed joint rfid estimation using snapshots of variable lengths," *Proc. of ACM Mobihoc*, pp. 247–256, 2015.
- [19] Q. Xiao, B. Xiao, and S. Chen, "Differential Estimation in Dynamic RFID Systems," *Proc. of IEEE INFOCOM*, pp. 295 – 299, April 2013.
- [20] H. Liu, W. Gong, L. Chen, W. He, K. Liu, and Y. Liu, "Generic Composite Counting in RFID Systems," *Proc. of IEEE ICDCS*, pp. 597 – 606, 2014.
- [21] Min Chen, Jia Liu, Shigang Chen, and Qingjun Xiao, "Anonymous category-level joint tag estimation: Poster," *Proc. of ACM Mobihoc*, pp. 363–364, 2016.
- [22] M. Chen and S. Chen, "Etap: Enable lightweight anonymous rfid authentication with o(1) overhead," *Proc. of IEEE ICNP*, pp. 267–278, November 2015.
- [23] A. Juels and S. A. Weis, "Defining Strong Privacy for RFID," *IEEE PerCom Workshops*, pp. 342 – 347, March 2007.
- [24] L. Lu, J. Han, R. Xiao, and Y. Liu, "ACTION: Breaking the Privacy Barrier for RFID Systems," *Proc. of IEEE INFOCOM*, 2009.
- [25] C. Tan, L. Xie, and Q. Li, "Privacy Protection for RFID-based Tracking Systems," *IEEE RFID*, 2010.
- [26] L. Lu, Y. Liu, and X. Li, "Refresh: Weak Privacy Model for RFID Systems," *Proc. of IEEE INFOCOM*, 2010.
- [27] T. Li, W. Luo, Z. Mo, and S. Chen, "Privacy-preserving RFID Authentication based on Cryptographical Encoding," *Proc. of IEEE INFOCOM*, 2012.
- [28] C. T. Nguyen, K. Hayashi, M. Kaneko, P. Popovski, and H. Sakai, "Probabilistic Dynamic Framed Slotted ALOHA for RFID Tag Identification," *Wireless Personal Communications*, vol. 71, pp. 2947–2963, August 2013.
- [29] S. Lee, S. Joo, and C. Lee, "An Enhanced Dynamic Framed Slotted ALOHA Algorithm for RFID Tag Identification," *Proc. of IEEE MobiQuitous*, 2005.
- [30] K. Y. Whang, B. T. Vander-Zanden, and H. M. Taylor, "A Linear-time Probabilistic Counting Algorithm for Database Applications," *ACM Transactions on Database Systems (TODS)*, vol. 15, no. 2, 1990.
- [31] Q. Chen, H. Ngan, Y. Liu, and L. M. Ni, "Cardinality Estimation for Large-Scale RFID Systems," *IEEE Transactions on Parallel and Distributed Systems*, vol. 22, no. 9, pp. 1441 – 1454, 2011.
- [32] N. L. Johnson and S. Kotz, *Urn Models and Their Application: An Approach to Modern Discrete Probability Theory*, John Wiley and Sons Inc, 1977.
- [33] *Stirling Numbers of the Second Kind*, Available at http://en.wikipedia.org/wiki/Stirling_numbers_of_the_second_kind.
- [34] *Zipf's Law*, Available at http://en.wikipedia.org/wiki/Zipf%27s_law.