# Digital Fingerprinting Codes: Problem Statements, Constructions, Identification of Traitors

Alexander Barg, *Senior Member, IEEE*, G. R. Blakley, and Grigory A. Kabatiansky

*Abstract*—We consider a general fingerprinting problem of digital data under which coalitions of users can alter or erase some bits in their copies in order to create an illegal copy. Each user is assigned a fingerprint which is a word in a fingerprinting code of size $M$ (the total number of users) and length $n$. We present binary fingerprinting codes secure against size-$t$ coalitions which enable the distributor (decoder) to recover at least one of the users from the coalition with probability of error $\exp(-\Omega(n))$ for $M = \exp(\Omega(n))$. This is an improvement over the best known schemes that provide the error probability no better than $\exp(-\Omega(n^{1/2}))$ and for this probability support at most $\exp(O(n^{1/2}))$ users. The construction complexity of codes is polynomial in $n$. We also present versions of these constructions that afford identification algorithms of complexity $\mathrm{poly}(n) = \mathrm{polylog}(M)$, improving over the best previously known complexity of $\Omega(M)$. For the case $t = 2$, we construct codes of exponential size with even stronger performance, namely, for which the distributor can either recover both users from the coalition with probability $1 - \exp(\Omega(n))$, or identify one traitor with probability 1.

*Index Terms*—Concatenated codes, fingerprinting problem, identification error, list decoding, polynomial-time decoding, separating codes.

## I. INTRODUCTION

LET $\mathcal{Q}$ be a finite alphabet of size $Q$. Suppose a dealer $D$ distributes copies of a long string $\Sigma$ over $\mathcal{Q}$ available by subscription to registered users of the system. Let $M$ be the total number of registered users. A copy assigned to the $i$th user contains a substring $\boldsymbol{x}(i) = (x_1(i), \ldots, x_n(i))$ (a fingerprint); fingerprints of different users are different. The goal of inserting the fingerprint is to personalize the copy given out to the user, and to rule out redistribution. Clearly, an individual user cannot

resell his copy of $\Sigma$ without running the risk of being tracked down. However, several users may collude in order to produce an unregistered copy. In doing so, they face the problem of supplying it with a fingerprint, which should, besides being different from each of their fingerprints, prevent $D$ from identifying members of the colluding group. Thus. the distributor $D$ faces the problem of constructing a large set of fingerprints (a fingerprinting code) that enables him to locate at least one member of the colluding group. The algorithm of creating registered fingerprints is parameterized by a secret key $k$. We assume that the algorithm is publicly known; however, the particular value of $k$ is kept secret by the distributor.

A group $U = \{u_1, \ldots, u_t\}$ of $t$ registered users that intend to produce an illegal copy of $\Sigma$ will be called a *coalition*. The goal of the coalition $U$ is to create a fingerprint $\boldsymbol{y} \in \mathcal{Q}^n$ of the illegal copy so that $D$ is unable to identify users from $U$. Following [5], we assume that the members of $U$ can only alter those coordinates of the fingerprint in which at least two of their fingerprints differ, and refer to this as the *Marking Assumption* (for more detailed discussion and motivation see [5]). Thus, it is known *a priori* that $y_m = x_m(u)$ for every $u \in U$ unless there is a pair $u, u' \in U$ such that $x_m(u) \neq x_m(u')$. What the users are allowed to do in the last case is a part of the formal description of the problem, made more precise in Section II-A.

The distributor faces the task of identifying one or more members of the coalition $U$ of traitors provided that $|U| = t$, where $t$ is a parameter. The fingerprinting problem thus is to design a set $C = \{\boldsymbol{c}(1), \ldots, \boldsymbol{c}(M)\}$ (a $t$-fingerprinting code, or more precisely, an ensemble of codes) in such a way that no matter which coalition $U$ of at most $t$ users collude to produce an unregistered fingerprint and no matter which algorithm the members of $U$ use, the distributor is always capable of tracing at least one of its members. This problem has been studied in different versions [1], [4]–[8], [13], [24], see especially [5] and [7] for an expanded informal discussion of the problem and literature overview. As customary in information and coding theory, we will address the existence question of families of $t$-fingerprinting codes with the number of codewords growing exponentially with the length of the code. In other words, to accommodate $M$ users, the distributor $D$ inserts order $\log M$ fingerprint digits. The answer depends on the rules that the members of $U$ (parents) are allowed to follow to create $\boldsymbol{y}$ (a descendant). As already remarked, when all the values $x_m(u)$ coincide, then $y_m$ must be of the same value. If some $x_m(u)$ differ, then the following two basic strategies employed by the users appear in the literature. Under the narrow-case fingerprinting problem, every coordinate of $\boldsymbol{y}$ can be chosen from the set of entries that its parents have in this coordinate [1], [6]–[8], [13]. Under the general

or wide-case problem, once there is a choice for $y_m$, it can be any letter from the alphabet. As suggested by [5], both cases can be expanded by allowing the users to make some symbols unreadable, or erased.

The narrow-case fingerprinting problem was expressly formulated in [13], where it was also proved that for $t = 2$ and $Q \geq 3$, there exist codes providing exact identification of at least one traitor with exponentially many codewords. For arbitrary $t \leq Q - 1$ this was proved in [1].

In the binary case, it turns out that exact identification of even one traitor is generally impossible. Therefore, we will allow some error rate, i.e., with low probability $\varepsilon$ the decision taken by the distributor can be wrong. Best known results in this problem were accomplished in [5], where the tradeoff between the size of a fingerprinting code and the probability $\varepsilon$ was considered. The best decline rate of the probability $\varepsilon$ accomplished in [5] is $\varepsilon = \exp(-\sqrt{n})$, irrespective of the code size (e.g., even for very small codes, for instance, of constant size). The question of attaining exponential code size and exponential decline of error rate with the length of the code $n$, natural from information-theoretic point of view, was left open in [5].

Furthermore, the construction of [5] in the part of identifying a guilty user relies upon complete (maximum-likelihood) decoding of a random code. The best known algorithms for this (NP-hard) problem require complexity of order $\Omega(M)$. In this paper, we will, relying on coding theory methods, present code families of size $\exp(\Omega(n))$ with $\varepsilon = \exp(-\Omega(n))$. Both the complexity of constructing codes and of their decoding grow as $\mathsf{poly}(\log M)$.

In Section II, we will give a precise statement of the problem with the aim of placing it in the standard information-theoretic context. This is worthwhile to do because rigorous statements of the digital fingerprinting problem do not appear in the literature, and moreover, because this enables us to establish relations between different versions of the problem. In formulating the problem, one faces a number of choices related to the worst versus average case. While in information theory there is an accepted standard, in the field of digital fingerprinting the situation is unsettled. We give a version of the definition geared toward worst case performance, which is consistent with the code constructions known in the literature and suggested in what follows, and also corresponds well to the nature of cryptographic problems. We conclude Section II by discussing previous results and goals of our paper.

One of the new ideas in code construction introduced in this work is the use of separating codes, which we briefly discuss in Section III, explaining also the reason for them to enter the fingerprinting problem. Section IV is devoted to constructions of binary fingerprinting codes. The constructions employ concatenation of two codes which proved useful in this problem (see [5]): a long outer code $W$ and a shorter inner code $V$. The code $W$ is error correcting with large minimum distance; the code $V$ has a $t$-separating property (see [14] and the survey [21]) and is used to separate law-abiding users from those committing fraud. The distance properties of $W$ enable us to amplify this separation. Section V gives identification algorithms of the distributor of complexity polynomial in the length of the codes constructed. This result is possible due to a surprising recent discovery of de-

coding algorithms of algebraic-geometry ("evaluation") codes that in polynomial time correct far more errors than half the minimum distance of the code [11].

## II. STATEMENT OF THE PROBLEM AND PREVIOUS RESULTS

### A. Problem Statements

Let us assume some ordering of the users and write them as $\{1, \ldots, M\}$. Let $\mathcal{Q}^n$ be the set of all $Q$-ary words of length $n$. For two words $\boldsymbol{x}, \boldsymbol{y} \in \mathcal{Q}^n$ let $\mathrm{d}(\boldsymbol{x}, \boldsymbol{y})$ denote the Hamming distance between them. Any subset $C \subset \mathcal{Q}^n$ is called a code of length $n$; if $|C| = M$ we say that $C$ is a code of size $M$ and denote its parameters by $(n, M)$. The number $R = R(C) := n^{-1} \log_Q M$ is called the *rate* of $C$. If $Q$ is a prime power and $C$ is a linear $\ell$-dimensional subspace of $\mathcal{Q}^n$, it is called a linear $[n, \ell, d]$ code. Here, $d$ is the minimum distance of $C$ (sometimes we omit it from the notation). The distance between a word $\boldsymbol{x}$ and a subset $\mathcal{H} \subset \mathcal{Q}^n$ is defined as $\mathrm{d}(\boldsymbol{x}, \mathcal{H}) := \min_{\boldsymbol{h} \in \mathcal{H}} \mathrm{d}(\boldsymbol{x}, \boldsymbol{h})$.

The distributor $D$ assigns to the user $i$ a fingerprint $\boldsymbol{x}(i) \in \mathcal{Q}^n$. The set of $M$ fingerprints $\{\boldsymbol{x}(1), \ldots, \boldsymbol{x}(M)\}$ forms a code $C \subset \mathcal{Q}^n$. As we show later (Propositions 2.3 and 2.6), in many important cases using a single code does not enable the distributor to solve the fingerprinting problem, namely, to locate a member of the coalition. Therefore, the distributor $D$ uses a family of codes $\mathcal{C} = \{C_k\}$, choosing a particular code $C_k$ with probability $\pi(k)$. More specifically, by the code $C_k$ we mean an *ordered* set of $M$ vectors $\{\boldsymbol{x}(1), \ldots, \boldsymbol{x}(M)\}$, so that $\pi(k) = P_M(\boldsymbol{x}(1), \ldots, \boldsymbol{x}(M))$. Here $k$ ranges over some set of $\mathcal{K}$ possible keys. In the constructions that follow, we will assume that $\pi(k) = \mathcal{K}^{-1}$. However, in the general problem statement of this section $\pi(k)$ can be an arbitrary distribution.

Having observed a fingerprint $\boldsymbol{y}$, the distributor $D$ identifies user $u$ as delinquent with some probability that depends both on $\boldsymbol{y}$ and on the specific code (key) used. Hence, the most general decision rule of $D$ can be described by a conditional probability distribution $P_D(u|\boldsymbol{y}, k)$. The distributions $\pi(k)$ and $P_D(u|\boldsymbol{y}, k)$ and the family of codes $\mathcal{C}$ are publicly known. The only information kept secret by $D$ is the specific value of $k$.

Let $U = \{u_1, \ldots, u_t\}$ be a coalition of $t$ users that collude to create an unregistered fingerprint $\boldsymbol{y}$. Let $X = C_k(U) = (\boldsymbol{x}^1, \ldots, \boldsymbol{x}^t)$ be the fingerprints assigned to the members of $U$ (for brevity we write $\boldsymbol{x}^i$ instead of $\boldsymbol{x}(u_i)$). The word $\boldsymbol{y}$ is taken from a subset $\mathcal{E}(X) \subset \mathcal{Q}^n$, called the *envelope* of $X$.

Particular ways to form the envelope are discussed later in more detail; for the moment, it can be assumed arbitrary. We assume that the members of $U$ attempt to confuse the distributor by choosing a particular value of $\boldsymbol{y}$ with some conditional probability $P_U(\boldsymbol{y}|X)$, where $X$ is the set of their fingerprints. Again, the distribution $P_U(\boldsymbol{y}|X)$ describes in the most general way the strategy that the members of $U$ use to confuse the distributor.

The distributor's probability of error in identifying a member of $U$ can be written in two ways

$$p_e(U, P_U)$$

$$= \sum_{u \notin U} \left[ \sum_{k=1}^{\mathcal{K}} \pi(k) \sum_{\boldsymbol{y} \in \mathcal{E}(C_k(U))} P_D(u|\boldsymbol{y}, k) P_U(\boldsymbol{y}|C_k(U)) \right]$$

(total probability of identifying user $u$ as guilty for all $u \notin U$), or

$$p_e(U, P_U)$$
$$= \sum_{k=1}^{\mathcal{K}} \pi(k) \left[ \sum_{u \notin U} \sum_{\boldsymbol{y} \in \mathcal{E}(C_k(U))} P_D(u|\boldsymbol{y}, k) P_U(\boldsymbol{y}|C_k(U)) \right]$$

(probability of error for a given code averaged over the choice of the code).

Let

$$p_e(U) = \max_{P_U} p_e(U, P_U)$$

be the error probability under the *optimal* strategy $P_U$ of the coalition (note that [5, Definition IV.2] implies a weaker attack, namely, only the uniform distribution on $\mathcal{E}(C_k(U))$. Thus, the maximum probability for $D$ to incriminate an user that is not a member of $U$ equals

$$p_e = p_e(\mathcal{C}, \pi(\cdot), P_D(\cdot)) = \max_{U, |U|=t} p_e(U).$$

This quantity still depends on the strategy $P_D(u|\boldsymbol{y}, k)$ of $D$, the choice of the family of codes $\{C_k\}$, and the probability $\pi(k)$ to choose a particular code $C_k$. It is natural to assume that the distributor performs these choices to minimize the value of $p_e$. The resulting probability is a function of the quadruple $(n, M, t; \mathcal{K})$ and the type of $\mathcal{E}$, and equals

$$\varepsilon(\mathcal{E}; n, M, t; \mathcal{K}) = \min_{\mathcal{C}, \pi(\cdot), P_D(\cdot)} p_e(\mathcal{C}, \pi(\cdot), P_D(\cdot)). \quad (1)$$

If the number of keys (codes in the family $\mathcal{C}$) is unrestricted, then we can introduce the value

$$\varepsilon(\mathcal{E}; n, M, t) = \min_{\mathcal{K}} \varepsilon(\mathcal{E}; n, M, t; \mathcal{K}). \quad (2)$$

Note that the total number of codes of size $M$ is bounded above by $2^{nM}$, which is also an upper bound on the number $\mathcal{K}$ of keys in this formula.

The general *fingerprinting problem* is to find $\varepsilon(\mathcal{E}; n, M, t)$ and codes and the optimal strategy of $D$ that achieve this error rate. With some abuse of language, a family of codes that enables $D$ to identify at least one member of $X$ as long as $|X| \le t$ (possibly, with error probability $\varepsilon$) will be called a *t-secure code*.

Given a $t$-subset $X \subset \mathcal{Q}^n$ we now define the envelope $\mathcal{E}(X)$. If $\boldsymbol{y} \in \mathcal{E}(X)$ then $\boldsymbol{y}$ is called a *descendant* of $X$ and any $\boldsymbol{x} \in X$ is called a *parent* of $\boldsymbol{y}$. Following [5], position $i$ is called *undetectable* for $X$ if the values of the words of $X$ match in their $i$th position: $x_i^1 = x_i^2 = \cdots = x_i^t$.

Denote by $\mathcal{Z}(X)$ the set of positions undetectable for $X$. By the marking assumption, the coalition cannot change the values of undetectable positions. If the position is detectable, then there are several options for the coalition to fill it. We will consider the narrow-sense and wide-sense envelopes and their expanded versions.

The *narrow-sense* envelope $e(X)$ is defined as follows:

$$e(X) = \{\boldsymbol{y} \in \mathcal{Q}^n | y_i \in \{x_i^1, \ldots, x_i^t\}\}.$$

The fingerprinting problem thus defined was studied in [1], [6], [7], [13], [24] (in the case of zero-error code constructions for the narrow case are also called codes with the identifiable parent

property or IPP codes). The *wide-sense* envelope $E(X)$ is defined as follows:

$$E(X) = \{\boldsymbol{y} \in \mathcal{Q}^n | y_i = x_i^1 \text{ for } i \in \mathcal{Z}(X)\}.$$

Both $e(X)$ and $E(X)$ can be generalized to the expanded case under which the coalition is allowed to generate unreadable (or erased) symbol $*$ in detectable positions. In particular, the expanded wide-sense envelope $E^*(X)$ is defined as follows (see [5, Definition II.3]):

$$E^*(X) = \{\boldsymbol{y} \in \{\mathcal{Q} \cup *\}^n | y_i = x_i^1 \text{ for } i \in \mathcal{Z}(X)\}.$$

It is clear that

$$\varepsilon(e; n, M, t; \mathcal{K}) \le \varepsilon(E; n, M, t; \mathcal{K}) \le \varepsilon(E^*; n, M, t; \mathcal{K}) \quad (3)$$

and

$$\varepsilon(e; n, M, t; \mathcal{K}) \le \varepsilon(e^*; n, M, t; \mathcal{K}) \le \varepsilon(E^*; n, M, t; \mathcal{K}). \quad (4)$$

Note that these inequalities are also valid for the probabilities $\varepsilon(\cdot; n, M, t)$ from (2).

### B. Properties of Different Versions of the Problem

In this subsection, we establish some relations between fingerprinting problems formulated above. We begin with a result that shows equivalence of the expanded wide-sense problem ($\mathcal{E} = E^*$) and the wide-sense one ($\mathcal{E} = E$).

*Proposition 2.1:* For any family of codes $\mathcal{C} = \{C_k\}$ and any probability distribution $\pi(k)$, the error rate of the distributor $\min_{P_D} p_e(\mathcal{C}, \pi(k), P_D)$ optimized over its strategies $P_D(\cdot)$ is the same for the expanded fingerprinting problem and the wide-sense one.

*Proof:* Clearly, the error rate achievable for the nonexpanded problem does not exceed the error rate for the expanded case. To show the converse, consider an arbitrary family of codes $\mathcal{C}$ for the $E$-case together with some probability distribution $\pi(k)$ defined on it. Let $P_D(u|\boldsymbol{y}, k)$ be the optimal decision rule for these codes $\{(C_k)\}$ and this $\pi(k)$. Then, for the expanded case, define the decision rule $P_D^*(\cdot)$ as follows: $P_D^*(u|\boldsymbol{y'}, k) = P_D(u|\boldsymbol{y}, k)$, where $y_i = \alpha$ if $y_i' = *$ and $y_i = y_i'$ otherwise, and where $\alpha$ is a fixed element of the alphabet. The transformation $\tau: \boldsymbol{y'} \to \boldsymbol{y}$ establishes a mapping from the set of strategies $P_U^*(\boldsymbol{y'}|C_k(U))$ of the coalition $U$ for the expanded wide-sense envelope $E^*$ on the set of strategies for the wide-sense envelope $E$, defined as follows:

$$P_U(\boldsymbol{y}|C_k(U)) = \sum_{\boldsymbol{y'} \in \tau^{-1}(\boldsymbol{y})} P_U^*(\boldsymbol{y'}|C_k(U)).$$

Clearly, the error probability $p_e$ of $D$ for the expanded case under the strategy $P_D^*(\cdot)$ and any given strategy of $U$ equals the error probability of $D$ for the nonexpanded case under the corresponding strategy of $U$. This proves the proposition. $\square$

The next corollary follows on replacing an arbitrary code family and a probability distribution by the optimal ones for the $E$-case.

*Corollary 2.2:* $\varepsilon(E; n, M, t; \mathcal{K}) = \varepsilon(E^*; n, M, t; \mathcal{K})$.

Note that for the narrow-sense envelope an analogous claim does not hold since, on the one hand, for the nonexpanded case there exist $t$-secure zero-error codes of exponential size (for $t < Q$), and on the other hand, for the expanded case the size of any $t$-secure code is at most $t$ (Proposition 2.6). The argument of the proof of Proposition 2.1 fails in this case at the step of setting up the correspondence $\boldsymbol{y}' \to \boldsymbol{y}$. Nevertheless, in the binary case, the one most often encountered in the literature, the wide-sense envelope is the same as the narrow-sense one. This leads to an important conclusion: *all the four statements of the general problem in the binary case coincide.*

In the problem of digital fingerprinting, special attention was devoted to the case of exact recovery ($\varepsilon = 0$). Note that the zero-error requirement is a very stringent condition. Indeed, if a single code $C$ is used and $t \geq Q$, then $|C| \leq Q$ even for $\mathcal{E} = e$, so zero-error $t$-secure, or $t$-IPP codes of nonzero rate do not exist. If $t \leq Q - 1$, it is possible to construct zero-error $t$-secure codes of exponential size if $\mathcal{E} = e$ (see [13] for the case $t = 2$ and [1] for arbitrary $t$).

However, nontrivial zero-error codes do not exist for the three other types of envelopes considered. Namely, we will establish that for the wide-sense envelope $E$ and for the expanded narrow-sense envelope $e^*$ not only the exact recovery $\varepsilon = 0$, but also relatively small error probability $\varepsilon$ cannot be achieved by a single code.

For a fixed-code family $\mathcal{C}$ we will write $p_e(\mathcal{C})$ to refer to error probability of the distributor in identifying a guilty user under the minimax formulation

$$p_e(\mathcal{C}) = \min_{\pi(k)} \min_{P_D} \max_{U \subset C, |U|=t} p_e(U).$$

We begin with a simple technical remark that for any a wide-sense $t$-secure code $C$ and its subcode $\mathcal{C} \subset C$

$$p_e(C) \geq p_e(\mathcal{C}). \tag{5}$$

Since the family consists of a single code, there is no distinction between a coalition and a subset of codewords, and the proof of (5) is obvious from the following inequality:

$$\begin{aligned} p_e(C) &= \min_{P_D} \max_{U \subset C, |U|=t} p_e(U) \\ &\geq \min_{P_D} \max_{U \subset \mathcal{C}, |U|=t} p_e(U) = p_e(\mathcal{C}). \end{aligned}$$

The next proposition builds upon the fact that if the system is based on a single code, then the members of the coalition know not only their fingerprints but also the fingerprints of all the other users.

*Proposition 2.3:* Let $C$ be a wide-sense $t$-secure code of size $|C| \geq 2t - 1$. Then

$$p_e(C) \geq \frac{t-1}{2t-1}. \tag{6}$$

*Proof:* We again identify the users and the codewords assigned to them. Let $\{\boldsymbol{x}^1, \boldsymbol{x}^2, \ldots, \boldsymbol{x}^t\}$ be the set of fingerprints assigned to a coalition $U$. The coalition $U$ creates a fingerprint $\boldsymbol{y}$ using the following deterministic strategy. Since the code $C$ is publicly known, $U$ chooses an arbitrary subcode $\mathcal{C} \subset C$ of size $2t - 1$ formed of its fingerprints and some other codewords $\boldsymbol{x}^{t+1}, \ldots, \boldsymbol{x}^{2t-1}$. Let $i$ be a given coordinate. If there exist at least $t$ vectors among $\{\boldsymbol{x}^1, \boldsymbol{x}^2, \ldots, \boldsymbol{x}^{2t-1}\}$ whose $i$th coordinates are all the same (and equal to some $a$) then the coalition $U$ necessarily contains one of these vectors. In this case, $U$ sets $y_i = a$. If this condition does not hold, then the $i$th position is detectable for the coalition $U$ and $y_i$ is set to $\alpha$, where $\alpha$ is some fixed element of the alphabet. It is easy to see that $\boldsymbol{y} = (y_1, \ldots, y_n) \in E(U)$ for any $U \subset \mathcal{C}, |U| = t$.

Consider any strategy of $D$, i.e., a set of probabilities $P_D(u|\boldsymbol{y})$. Suppose, without loss of generality, that

$$P_D(\boldsymbol{x}^1|\boldsymbol{y}) \leq P_D(\boldsymbol{x}^2|\boldsymbol{y}) \leq \cdots \leq P_D(\boldsymbol{x}^{2t-1}|\boldsymbol{y}).$$

Then, the probability of identifying a member of the coalition $U = (\boldsymbol{x}^1, \ldots, \boldsymbol{x}^t)$ is

$$P_D(\boldsymbol{x}^1|\boldsymbol{y}) + \cdots + P_D(\boldsymbol{x}^t|\boldsymbol{y}) \leq \frac{t}{2t-1}.$$

Hence,

$$p_e(C) \geq 1 - \frac{t}{2t-1} = \frac{t-1}{2t-1}.$$

Conclude by (5). $\qquad\square$

*Corollary 2.4:* Let $M \geq 2t - 1$. Then

$$\varepsilon(E; n, M, t; \mathcal{K}) \geq \frac{t-1}{\mathcal{K}(2t-1)}.$$

*Proof:* Let $\mathcal{C} = \{C_k, 1 \leq k \leq \mathcal{K}\}$ be a family of wide-sense $t$-secure codes. There is a key $k_0$ such that its probability $\pi(k_0) \geq \mathcal{K}^{-1}$. Consider the same strategy of the coalition as in the proof of Proposition 2.3, assuming that $k = k_0$. Then

$$P_e \geq \pi(k_0)P_e(C_{k_0}) \geq \frac{t-1}{(2t-1)\mathcal{K}}. \qquad\square$$

In what follows, we do not consider codes with $\leq t$ codewords, which we call trivial. The proof of Proposition 2.3 also shows that for any code $C$ such that $|C| = t + j$, where $1 \leq j \leq t - 1$, the error rate $p_e(C) \geq j/(t+j)$. Hence, the error rate $p_e(C) \geq 1/(t+1)$ for any nontrivial code $C$. We obtain the following result.

*Corollary 2.5:* Let $C$ be a wide-sense $t$-secure zero-error code. Then $|C| \leq t$.

Recall a result in [5] which shows that zero-error recovery is impossible with a single binary $t$-secure code $C$ with $|C| \geq 3$. The above sequence of results develops this by establishing a limiting tradeoff between the error rate and the size of the family for any alphabet size $Q$; in particular, for a single code we have the estimate (6).

The same reasoning can be applied to the expanded narrow-sense problem.

*Proposition 2.6:* The results of Proposition 2.3 and Corollaries 2.5 and 2.4 are valid if the wide-sense problem is replaced by the expanded narrow-sense problem.

*Proof:* Let $\mathcal{C} = \{\boldsymbol{x}^1, \ldots, \boldsymbol{x}^{2t-1}\}$ be a subcode of the code $C$. Let us modify the strategy of the coalition defined in the proof of Proposition 2.6, replacing $\alpha$ with $*$ and leaving the rest unchanged. Now the argument of that proof can be applied verbatim to the case considered. $\qquad\square$

Thus, we arrive at the following important conclusion. Under all problem statements except the narrow case, we have

$$\mathcal{K} \geq \frac{1}{(t+1)\varepsilon}.$$

In other words: *to achieve exponential decline of the error rate the number of keys has to grow exponentially with the length of the fingerprint*. Note that this growth order of the number of keys will be achieved in the constructions of Sections IV and V.

### C. Previous Results

We have seen that in the case of binary codes, first considered in Boneh and Shaw [5], all the versions of the fingerprinting problem considered above coincide. It is proved in [5] that it is not possible to solve the fingerprinting problem by using a fixed assignment of fingerprints to the users and suggested instead to use random choice of a code from a given code family. This enabled the authors of [5] to construct families of $t$-secure binary ($Q = 2$) fingerprinting codes. This construction is the best known in the literature from the point of view of code parameters.

*Theorem 2.7 [5, Theorem V.5]:* There exists a family of $t$-secure binary fingerprinting $(n, M)$ codes with error probability $\varepsilon$, where

$$n = O(t^4 (\log \varepsilon / M) \log \varepsilon).$$

Solving this for $\varepsilon$ and putting $\log M = O(n^\alpha)$, we obtain $\log \varepsilon = -\Omega(\min(\sqrt{n}, n^{1-\alpha}))$. Thus, in the construction of [5], $\varepsilon$ cannot decrease faster than $\exp(-\Omega(\sqrt{n}))$ and then $M = 2^{O(\sqrt{n})}$. On the other hand, as proved in [5, Theorem VI.1], for any family of codes

$$\varepsilon \geq \frac{1}{t} \exp\left(-\frac{2}{t-3} n\right)$$

and so it might be possible to construct codes with error probability falling exponentially with $n$. For this reason, [5] puts forward the question of tightening the gap between the known bounds and constructions. The main goal of this paper is to resolve this question by constructing codes with exponentially small probability of identification error and the number of words growing exponentially with $n$, and to show that this performance is attainable with polynomial complexity of code construction and decoding.

Another parameter of the problem whose importance has not been singled out in previous works, but was emphasized in this section, is the number of keys in the fingerprinting scheme. It is desirable to have as few keys as possible because the distributor must store the key bits in order to manage the system. The best known scheme [5] does not invoke this parameter explicitly, referring instead to a random code of length $\Omega(n)$ and size $M$. Suppose that $M = \exp(\Omega(n))$, then the total number of keys is $\mathcal{K} = \Omega(2^{nM})$, i.e., $\log \mathcal{K} = \exp(\Omega(n))$, and hence $D$ must store $\exp(\Omega(n))$ bits of information. In contrast, we will obtain $\log \mathcal{K} = \Omega(n)$, which by Proposition 2.6 is the best possible order of magnitude.

## III. SEPARATING CODES

In this section, we establish relationships between separating property of codes and fingerprinting. A $Q$-ary code $V$ is called $(t, t')$-separating [14] if for any two disjoint subsets $X \subset V$, $Y \subset V$ such that $|X| = t$, $|Y| = t'$ there holds

$$e(X) \cap e(Y) = \emptyset.$$

In other words, for any two disjoint subsets $X$, $Y$ with $|X| = t$, $|Y| = t'$ there is a coordinate $i$ such that

$$\{x_i, \boldsymbol{x} \in X\} \cap \{y_i, \boldsymbol{y} \in Y\} = \emptyset.$$

In what follows, we assume that $t + t' \geq 3$ because the case $t = t' = 1$ is trivial.

Separating codes were studied in [14], [18]–[21], [15], and also under the name of secure frameproof codes in [25] and of partially identifying codes in [8].

*Lemma 3.1:* For any single code $V$ which is not $(t, t)$-separating, the identification error probability $p_e(V) \geq 1/2$ for the narrow-sense as well as the wide-sense envelope, expanded or not.

*Proof:* If the $(t, t)$-separating property does not hold, then there exist two disjoint subsets $U, U' \subset V$ such that $|U| = t$, $|U'| = t$, and $e(U) \cap e(U') \neq \emptyset$. Let $\boldsymbol{y} \in e(U) \cap e(U')$. Suppose that the strategy of both coalitions $U$ and $U'$ consists of choosing $\boldsymbol{y}$ as the fingerprint. Then, irrespective of the strategy of $D$, the error probability of deciding between $U$ and $U'$, and thus the worst case error $p_e(V)$ is bounded below by $1/2$. $\qquad\square$

It was noted in the literature that separation is necessary (but not sufficient) for narrow-sense zero-error identification. In the context of probabilistic fingerprinting this remark is new.

*Lemma 3.2:* For any single $(t, t)$-separating code $V$, the identification error probability $p_e(V) \leq 1 - 1/t$ for the narrow-sense as well as for the wide-sense envelope, expanded or not.

*Proof:* Consider any size-$t$ coalition $U \in V$ and any strategy $P_U$ of generating the fingerprint $\boldsymbol{y}$. The distributor's probability of correctly identifying a member of $U$ is at least $1/t$ by the following strategy: $D$ finds any coalition $U'$ such that $\boldsymbol{y} \in E(V(U'))$ and identifies randomly a user $u \in U'$. Since $V$ is $(t, t)$-separating, we have $|U' \cap U| \geq 1$; thus, $u \in U$ with probability $\geq 1/t$. $\qquad\square$

We are especially interested in the particular case $t = 2$, which received extensive coverage in the literature (see the survey [21]). In this case, the identification error probability of any single $(2, 2)$-separating $(n, M)$ code $V$ equals $p_e(V) = 1/3$ for narrow- or wide-sense envelope, expanded or not. In three out of four cases, this is optimal for a single fingerprinting code.

*Lemma 3.3:* Let $V$ be a $(2, 2)$-separating $(m, M)$ code. Then

$$p_e(V) = \varepsilon(e^*; m, M, 2; \mathcal{K}) = \varepsilon(E; m, M, 2; \mathcal{K})$$
$$= \varepsilon(E^*; m, M, 2; \mathcal{K}) = 1/3.$$

*Proof:* By Propositions 2.3 and 2.6, we have $p_e(V) \geq 1/3$ for $\mathcal{E} = e^*$, $E$ or $E^*$. On the other hand, let $\boldsymbol{y}$ be a fingerprint ob-

served by the distributor. Consider the set of all pairs $\boldsymbol{x}$, $\boldsymbol{x}'$ such that $\boldsymbol{y} \in E(\boldsymbol{x}, \boldsymbol{x}')$. By the $(2, 2)$-separating property, either all such pairs have a common element, or there are three pairs that form a triangle configuration $\{(\boldsymbol{x}, \boldsymbol{x}'), (\boldsymbol{x}', \boldsymbol{x}''), (\boldsymbol{x}'', \boldsymbol{x})\}$. In the first case, $D$ outputs the common element, and then $p_e(V) = 0$. In the last case, the strategy of $D$ is to take a random element from the triple $\boldsymbol{x}$, $\boldsymbol{x}'$, $\boldsymbol{x}''$, and then $p_e(V) = 1/3$. $\square$

Hereafter, we focus on the binary case. Let

$$R_{t, t'}^{(s)}(m) = \max_{\substack{V \subset \{0,1\}^m \\ V \text{ is } (t,t')\text{-separating}}} \frac{\log_2 |V|}{m}$$

be the maximum rate of $(t, t')$-separating codes of length $m$ and let

$$R_{t, t'}^{(s)} = \lim_{m \to \infty} R_{t, t'}^{(s)}(m) \qquad (7)$$

be the maximum achievable rate of $(t, t')$-separating codes (as usual, it is not known if this limit exists; if it does not, then the bounds below should be formulated under the $\liminf$ or $\limsup$ definition, as appropriate). If $t = t'$, we use a short notation $R_t^{(s)}$. Existence bounds on $(2, 2)$-separating codes were studied in [14], [18]; the best known bound [18] gives for the asymptotic rate of linear codes the value

$$R = \ell/m \geq (3 - \log_2 7)/3 \approx 0.0642.$$

On the other hand, it is known that asymptotically $R \leq 0.108$ for linear codes [21] and $R_2^{(s)} \leq 0.283$ for arbitrary codes (see [15], [21]). Paper [21] also contains tables of $(2, 2)$-separating $[m, \ell]$ codes. For instance, the well-known $[7, 3, 4]$ code is readily seen to be $(2, 2)$-separating. Below, we use a $[35, 6]$ code from [21] at the inner level to construct a 2-secure fingerprinting code with very low probability of identification error (see Example after Theorem 4.4).

It is clear that binary linear $(t, t')$-separating codes do not exist for $\max(t, t') > 2$. In the unrestricted case, the existence of $(t, t')$-separating codes is established by a standard probabilistic argument ("random coding with expurgation"); see [20], [15] for $t = t' = 2$. We have the following.

*Proposition 3.4:* Let $T = t + t'$. There exist binary $(t, t')$-separating codes of length $m$ and size $\frac{1}{2} P_T^{-\frac{m}{T-1}}$, where $P_T = 1 - 2^{-(T-1)}$, i.e.,

$$R_{t, t'}^{(s)}(m) \geq -\frac{\log_2 P_T}{T - 1} - \frac{1}{m}.$$

*Proof:* Consider a random binary $(m, L)$ code $V$ and compute the expectation $\boldsymbol{E}$ of the number of pairs of subsets $X, Y$ of the code $V$, $|X| = t$, $|Y| = t'$ that contradict the $(t, t')$ separating property. Whenever $\boldsymbol{E} \leq L/2$ then $(m, L/2)$ codes with the $(t, t')$ separating property exist.

The probability that a given pair $X$ and $Y$ violate the separating property is

$$\left( 1 - 2^{-(T-1)} \right)^m = P_T^m.$$

The expectation of the number of pairs $X, Y$ that violate the separating property is

$$\boldsymbol{E} \leq \binom{L}{t} \binom{L - t}{t'} P_T^m.$$

Take

$$L = \left( \frac{t! t'!}{-m} P_T^{-m} \right)^{1/(T-1)}$$

then

$$\boldsymbol{E} < \frac{L^T}{t! t'!} P_T^m = \frac{L}{2}.$$

Finally, note that $(t! t'!/2)^{(1/(T-1))} \geq 1$. $\square$

Note that [14] gives an asymptotic bound $R_{t, t'}^{(s)} \geq -\frac{\log_2 P_T}{T}$ which is somewhat weaker than this result.

## IV. CODE CONSTRUCTIONS

### A. The Case of Arbitrary $t$

We will be concerned with binary codes $(Q = 2)$. Recall again that all the statements of the fingerprinting problem in this case coincide. To construct a family of fingerprinting codes $\mathcal{C}$ we use the idea of concatenation [10]. Recall that a binary concatenated code of length $mN$ and size $M$ is formed by a binary inner $(m, q)$ code $V$ and an $q$-ary outer code $W$ of length $N$ and size $M$ based on a fixed one-to-one mapping of the $q$-ary alphabet to $V$. For every codeword $\boldsymbol{w} \in W$, the word of the concatenated code is obtained by replacing coordinates of $\boldsymbol{w}$ with the corresponding codewords of $V$.

Consider as the inner code $V$ a binary $(t, t)$-separating code of length $m$ and size $q$. We choose a $q$-ary $[N, K, \Delta = \delta N]$ linear code as an outer code $W$. Since we need to obtain a family of codes $\mathcal{C}$, our construction also involves $N$ random bijections $\varphi_i \colon \mathbb{F}_q \to V$ $(i = 1, 2, \ldots, N)$, where $\mathbb{F}_q$ is the field of $q$ elements. Restrictions on the parameters of the codes $V$ and $W$ will become clear when we analyze the construction.

Consider vector mappings $\phi = (\varphi_1, \varphi_2, \ldots, \varphi_N)$, where each $\varphi_i$, $1 \leq i \leq N$ is an arbitrary bijection. Let us number these mappings from 1 to $(q!)^N$ in an arbitrary order and write $\phi^{(k)}$ to refer to the $k$th mapping. A typical codeword of $C_k$ is obtained by taking a codeword $\boldsymbol{w} = (w_1, \ldots, w_N) \in W$ and computing the binary $mN$-vector

$$\phi^{(k)}(\boldsymbol{w}) = \left( \varphi_1^{(k)}(w_1), \varphi_2^{(k)}(w_2), \ldots, \varphi_N^{(k)}(w_N) \right).$$

The mapping $\phi^{(k)}$ is chosen by $D$ with uniform distribution and constitutes the secret key. The length of the code $C_k$ is $mN$, the size $M = q^K$. Each user is assigned a fingerprint given by a code vector in $C_k$. It will be convenient to identify the users with the code vectors of the code $W$ (note that $|W| = |C_k|$). The $(t, t)$-separating property of $V$ will be essential for the entire scheme.

Assume that a coalition $U = \{\boldsymbol{u}^1, \ldots, \boldsymbol{u}^t\} \subset W$ generates a fingerprint

$$\boldsymbol{y} = (\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_N) \in E\left( \varphi^{(k)}(\boldsymbol{u}^1), \ldots, \varphi^{(k)}(\boldsymbol{u}^t) \right)$$

where the subblocks $\boldsymbol{y}_i$ are length-$m$ binary vectors. The decision algorithm of the distributor $D$ first proceeds with "decoding" of every vector $\boldsymbol{y}_i$ with the inner code $V$ and then decodes the result with the outer code $W$.

Algorithm 1:

1) The first stage is decoding of the code $V$. For a given $i$ take an arbitrary $t$-set $\{\hat{\boldsymbol{v}}^1, \ldots, \hat{\boldsymbol{v}}^t\}$ of code vectors from $V$ that can generate $\boldsymbol{y}_i$, i.e., such that $\boldsymbol{y}_i \in E(\hat{\boldsymbol{v}}^1, \ldots, \hat{\boldsymbol{v}}^t)$. Define the set $H_i = \{h_i^1, \ldots, h_i^t\}$, where $\varphi_i^{(k)}(h_i^j) = \hat{\boldsymbol{v}}^j$. The result of the first-stage decoding is given by the set

$$\mathcal{H} = \{\boldsymbol{h} \in F_q^N : h_1 \in H_1, \ldots, h_n \in H_N\}.$$

2) The distributor $D$ identifies a possible member of the coalition $U$ as the user that corresponds to the code vector $\hat{\boldsymbol{w}} \in W$ such that

$$s(\hat{\boldsymbol{w}}, \mathcal{H}) = \max_{\boldsymbol{w} \in W} s(\boldsymbol{w}, \mathcal{H})$$

where

$$s(\boldsymbol{w}, \mathcal{H}) := |\{i : w_i \in H_i\}| = N - \mathrm{d}(\boldsymbol{w}, \mathcal{H}).$$

We note that this algorithm has two nonstandard features. The first-stage decoding outputs a subset (a *list*) of the words. This differs from the standard decoding algorithm of concatenated codes (which is also [5, Algorithm 2]), where it is assumed that this output is a single codeword. In the second stage, the algorithm looks for a nearest code vector of $W$ to the *subset* $\mathcal{H}$ rather than to one "received" word.

The reason that this algorithm and its subsequent modifications will work for the identification problem is that the number of occasions that coordinates of a code vector fall in the set $H_i$ is a random variable whose average is greater when the code vector is a member of the coalition $U$ than when it is not. Hence, the probability that the distributor identifies incorrectly a code vector as a member of the coalition falls exponentially with the length $N$ of the code $W$.

We will use the following standard bounds on large deviations. Let $\xi_i$ be independent Bernoulli random variable equal to 1 with probability $p$ and 0 with probability $1-p$. Then the probabilities of the tails can be bounded as

$$\Pr\left\{\sum_{i=1}^N \xi_i \leq N\sigma\right\} \leq 2^{-ND(\sigma\|p)}, \qquad \text{if } p > \sigma$$

$$\Pr\left\{\sum_{i=1}^N \xi_i \geq N\sigma\right\} \leq 2^{-ND(\sigma\|p)}, \qquad \text{if } p < \sigma$$

where $D(\sigma\|p) = \sigma \log_2(\sigma/p) + (1-\sigma)\log_2((1-\sigma)/(1-p))$ is the information divergence of two binomial distributions.

*Theorem 4.1:* Let $W$ be a $q$-ary linear $[N, K, \Delta = \delta N]$ code with

$$\delta > 1 - \frac{1}{t^2} + \frac{t-1}{t(q-1)}$$

and let $V$ be a $(t, t)$-separating binary $(m, q)$ code. The family of concatenated codes $\mathcal{C} = \{C_k\}$ with inner code $V$, outer code $W$, and the set of all $(q!)^N$ bijections $\phi^{(k)}$, together with the

decision rule defined by Algorithm 1, forms a binary fingerprinting code of length $n = mN$ with $q^K$ code vectors (users) that identifies one traitor with error probability

$$p_e \leq 2^{-nR(V)[(\log_2 q)^{-1} D(\sigma\|\frac{t-1}{q-1}) - R(W)]} \qquad (8)$$

where $\sigma = 1/t - (1-\delta)t$.

*Proof:* Let $U_i = \{u_i^1, \ldots, u_i^t\}$ be the multiset of the $i$th coordinates of the vectors in the coalition $U$. We denote the number of distinct elements in $U_i$ by $|U_i|$. Since for the sets $H_i$ and $U_i$, their images under the bijection $\varphi_i^{(k)}$ can generate the same vector $\boldsymbol{y}_i$ and since the inner code $V$ is $(t, t)$-separating, we have $H_i \cap U_i \neq \emptyset$. This implies that $\sum_{j=1}^t s(\boldsymbol{u}_j, \mathcal{H}) \geq N$, and hence $\max_j s(\boldsymbol{u}_j, \mathcal{H}) \geq N/t$.

On the other hand, for any $\boldsymbol{w} \notin U$, the element $w_i$ can be included in the list $H_i$ for one of the two reasons. The first possibility is that $w_i \in U_i$, and the number $T$ of such positions is at most $t(N - \Delta)$ since any two distinct vectors from $W$ coincide in at most $N - \Delta$ positions. The other option is that the mapping of $w_i$ to $V$ matches the random bijection $\varphi_i^{(k)}$. To compute the probability $p_i$ of the last event, note that, by assumption, we have $w_i \notin U_i$, thus, $w_i \notin (U_i \cap H_i)$. Since the bijections are chosen randomly and uniformly, we obtain

$$p_i := \Pr\{w_i \in H_i | w_i \notin U_i\} = \frac{t - |U_i \cap H_i|}{q - |U_i|} \leq \frac{t-1}{q-1}. \quad (9)$$

Now let us bound above the probability that $s(\boldsymbol{w}, \mathcal{H}) \geq N/t$ for a particular vector $\boldsymbol{w}$. Let $\xi_i$ be independent Bernoulli random variables equal to 1 with probability $p_i$ and 0 with probability $1 - p_i$. Then, for $\boldsymbol{w} \notin U$

$$\Pr\{s(\boldsymbol{w}, \mathcal{H}) \geq N/t | \boldsymbol{w} \notin U\}$$
$$\leq \Pr\left\{\sum_{i=1}^{N-T} \xi_i \geq \frac{N}{t} - T\right\}$$
$$\leq \Pr\left\{\sum_{i=1}^N \xi_i \geq \frac{N}{t} - t(N-\Delta)\right\} \leq 2^{-ND(\sigma\|\frac{t-1}{q-1})} \quad (10)$$

where in the last step we relied upon the inequality $\sigma > (t-1)/(q-1)$ implied by the condition on $\delta$ in the statement of the theorem. Now

$$p_e \leq \Pr\left\{\max_{\boldsymbol{w} \notin U} s(\boldsymbol{w}, \mathcal{H}) \geq N/t\right\}$$
$$\leq q^K \Pr\{s(\boldsymbol{w}, \mathcal{H}) \geq N/t | \boldsymbol{w} \notin U\}.$$

Let us substitute the bound on the last probability. Taking the logarithm of the right-hand side and using the definitions $R(W) = K/N$ and $R(V) = \log_2 q/m$, we obtain

$$p_e \leq 2^{KmR(V)} 2^{-ND(\sigma\|\frac{t-1}{q-1})}$$

which is the same as (8). □

*Remark:* The theorem yields nontrivial results, i.e., the exponential decline of the distributor's probability of incorrectly identifying a user as guilty, when the outer code $W$ has large distance $\Delta$. It is well known [19] (rediscovered in [25]) that concatenation of an inner separating code with an error-correcting outer code with large distance ($\delta > 1 - 1/t^2$) gives a longer separating code. We would like to stress that the separating property (renamed in crypto literature as secure frameproof property) is

not sufficient for identifying members of the coalition. Therefore, in the last theorem we perform a probabilistic analysis of the random ensemble of codes which amplifies the separating property to achieve a desired fingerprinting performance. Note also that concatenation of a fingerprinting inner code with an error-correcting code used in [5] does not lead to asymptotically good fingerprinting codes.

Let us use the general construction of this theorem together with specific choices of codes to present a few families of easily constructible $t$-secure fingerprinting codes. These constructions will enable us to claim overall code rate separated from zero together with exponentially small probability of identification error for $t = \text{const.}$ We will tacitly assume that for a fixed choice of the codes $V$ and $W$, the family of codes $\{C_k\}$ is obtained by applying $N$ random mappings $\varphi_i^{(k)} \colon \mathbb{F}_q \to V$, $i = 1, 2, \ldots, N$ to the coordinates of vectors in $W$. In the following, we focus only on the choice of the codes.

Take $V$ to be a long binary $(t, t)$-separating code of length $m$ and rate

$$R_{2t} = m^{-1} \log_2 q = -\frac{\log_2 P_{2t}}{2t - 1}$$

(see Proposition 3.4) and $W$ a $[q, K]$ extended Reed–Solomon (RS) code of rate $R(W) = K/q$ over $\mathbb{F}_q$. Substituting these parameters into Theorem 4.1 and taking into account that

$$(\log_2 q)^{-1} D\left(\sigma \left\| \frac{t-1}{q-1}\right.\right) \sim \sigma$$

for $t$ fixed and $q$ growing, we obtain the following result.

*Corollary 4.2:* For any fixed $t$ and any rate $R < R_{2t}/t(t+1)$, a binary fingerprinting code of length $n$ with $2^{Rn}$ code vectors (users) constructed by concatenating the codes $V$ and $W$ identifies one traitor with exponentially decreasing error probability

$$p_e \leq 2^{-n[t^{-1}R_{2t} - (t+1)R + o(1)]}.$$

Now let us take $V$ to be a fixed binary $(t, t)$-separating $(m, q)$ code, where $q$ is an even power of a prime. This choice is possible due to Proposition 3.4. Let us take the outer code $W$ to be a long algebraic-geometry (AG) code from a maximal curve, whose parameters asymptotically approach the bound $R(W) = 1 - \delta - 1/(\sqrt{q} - 1)$[27].

*Corollary 4.3:* For any fixed $t$, any $q \geq (t^2 + 2)^2$ that is an even power of a prime, and any rate $R < R_{2t}R_0(W)$, where $R_0(W)$ is the root of the equation

$$R(W) \log_2 q = D\left(\frac{1}{t} - t\left(R(W) - \frac{1}{\sqrt{q}-1}\right) \left\| \frac{t-1}{q-1}\right.\right) \tag{11}$$

a binary fingerprinting code of length $n$ and size $2^{Rn}$ constructed by concatenating a fixed inner $(m, q)$ code $V$ with the $(t, t)$-separating property and AG codes $W$ of rate $R(W) \leq R_0(W)$ and growing length $N$, identifies one traitor with exponentially falling error probability given by (8).

*Proof:* By Proposition 3.4, for any $m$, there exists a binary $(t, t)$-separating code $V$ of length $m$, size

$$|V| \geq (1/2)\left(1 - 2^{-(2t-1)}\right)^{-m/(2t-1)}$$

and, thus, of rate $R(V) \approx R_{2t}$. Hence, for a constant $t$ and sufficiently large $m$, it is possible to find a code of size $q$ such that $q$ satisfies the conditions of the corollary. Consider a $q$-ary AG code $W$ of length $N$ and rate $R(W) < R_0(W)$ such that for large $N$, its distance $\Delta = \delta N$ satisfies $1 - \delta \leq R(W) + 1/(\sqrt{q} - 1)$.

By Theorem 4.1, the error probability of identification will fall exponentially if $\sigma > (t - 1)/(q - 1)$ and

$$0 < R(W) \log_2 q < D\left(\sigma \left\| \frac{t-1}{q-1}\right.\right). \tag{12}$$

The inequality

$$\sigma = \frac{1}{t} - t\left(R(W) + \frac{1}{\sqrt{q}-1}\right) > \frac{t-1}{q-1}$$

is equivalent to the inequality $R(W) < A$, where

$$A = \frac{1}{t^2} - \frac{t-1}{t(q-1)} - \frac{1}{\sqrt{q}-1}.$$

It is immediate to verify that the condition $q \geq (t^2 + 2)^2$ implies that $A > 0$. Hence, the upper bound (8) holds true if $0 < R(W) < A$.

Next, let us show that the segment of values of $R(W)$ that satisfy (12) is a proper subset of $(0, A)$. This is immediate since for $R(W) \in (0, A)$, $D(\sigma \| \frac{t-1}{q-1})$ is a positive decreasing function of $R(W)$ which reaches zero when $\sigma = (t-1)/(q-1)$, and the middle part of (12) increases from zero to $A \log_2 q$. Hence, (11) has a (single) root, $R_0(W) < A$, and for all $0 < R(W) < R_0(W)$, inequality (12) holds true. This completes the proof. $\square$

Note that the codes of Corollary 4.3 are polynomially constructible. Indeed, by a recent result of [22], the construction complexity of the code $W$ for them is $O((N \log_q N)^3)$, and the complexity of constructing the code $V$ is constant, independent of $N$.

### B. The Case $t = 2$

A usual assumption in digital fingerprinting is that it is not possible to recover "the entire coalition since some of its members might be passive" (see, e.g., [5, p. 1899]). Therefore, the fingerprinting problem was restricted to finding one guilty user with high probability. In the case of size-2 coalitions, it is possible to construct a family of fingerprinting codes $\mathcal{C}$ that have a stronger property than in the original definition. Namely, the distributor either recovers one member of the coalition $U = \{\boldsymbol{u}, \boldsymbol{u}'\}$ with probability one (*zero-error*), or *both* members with probability $1 - \varepsilon$.

To construct a family of fingerprinting codes $\mathcal{C} = \{C_k\}$ in the case of $t = 2$, we use the same general idea of concatenation as for arbitrary $t$, with a somewhat different decision algorithm. The construction involves a binary $(2, 2)$-separating $[m, \ell]$ linear code $V$ [18], a $q$-ary $[N, K, \Delta = \delta N]$ code, $q = 2^\ell$, and $N$ random bijections $\varphi_i^{(k)} \colon \mathbb{F}_q \to V$, $i = 1, \ldots, N$.

Encoding, or the fingerprinting assignment procedure, is the same as above. Namely, the fingerprint corresponding to a vector $\boldsymbol{w} = (w_1, \ldots, w_N) \in W$ is obtained by computing $\boldsymbol{x} = (\varphi_1^{(k)}(w_1), \ldots, \varphi_N^{(k)}(w_N))$. The length of the code $C_k$ is $n = mN$ and the size is $q^K$. As before, we identify the users with the code vectors of the code $W$.

Assume that a coalition $U = \{\boldsymbol{u}, \boldsymbol{u}'\} \subset W$ generates a fingerprint $\boldsymbol{y} = (\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_N)$, where the subblocks $\boldsymbol{y}_i$ are length-$m$ binary vectors.

Let us first describe inner decoding, which will be different from the general case. Recall that $Q = 2$, and so the narrow-sense envelope $e$ and the wide-sense envelope $E$ are the same. By definition, the subblock $\boldsymbol{y}_i$ is contained in the envelope $e(\{\boldsymbol{v}, \boldsymbol{v}'\})$, where $\boldsymbol{v} = \varphi_i^{(k)}(u)$ and $\boldsymbol{v}' = \varphi_i^{(k)}(u')$. Consider all possible pairs $(\boldsymbol{v}^1, \boldsymbol{v}^2)$ of code vectors from the code $V$ such that $\boldsymbol{y}_i \in e(\{\boldsymbol{v}^1, \boldsymbol{v}^2\})$. By the $(2, 2)$-separating property, all these pairs should intersect. Hence, either there is a vector $\boldsymbol{v}^*$ that is an element of every such pair (a *star* configuration), or there are three such pairs that form a triangle $\{\boldsymbol{v}, \boldsymbol{v}', \boldsymbol{v}''\}$ (cf. Lemma 3.3). In the first case, the decoding result is given by $H_i = (\varphi_i^{(k)})^{-1}(\boldsymbol{v}^*) \in \mathbb{F}_q$; note that $H_i \in \{u_i, u_i'\}$. In the second case, the result is given by the *list*

$$H_i = \left\{ (\varphi_i^{(k)})^{-1}(\boldsymbol{v}), (\varphi_i^{(k)})^{-1}(\boldsymbol{v}'), (\varphi_i^{(k)})^{-1}(\boldsymbol{v}'') \right\}$$

and $\{u_i, u_i'\} \subset H_i$.

Algorithm 2:

1) The $N$ columns $\boldsymbol{y}_1, \ldots, \boldsymbol{y}_N$ of the observed fingerprint $\boldsymbol{y}$ are independently decoded with the code $V$. The decoding result $H_i$ of the $m$-vector $\boldsymbol{y}_i$ is either a single element of $\mathbb{F}_q$ or a 3-subset of $\mathbb{F}_q$. The result of this step is the set

$$\mathcal{H} = \{\boldsymbol{h} \in \mathbb{F}_q^N : h_1 \in H_1, \ldots, h_N \in H_N\}.$$

2) Denote by $\mathcal{I}$ the subset of coordinates in which $|H_i| = 1$ and let $\mathcal{L} = \{1, 2, \ldots, N\} \backslash \mathcal{I}; I = |\mathcal{I}|, L = |\mathcal{L}|$. If there is a vector $\boldsymbol{w} \in W$ such that the number of agreements

$$|\{i \in \mathcal{I} : w_i = H_i\}| > 2(N - \Delta)$$

then the distributor identifies $\boldsymbol{w}$ as a member of $U$.

3) If the condition of the previous step is not satisfied, the distributor finds all code vectors $\boldsymbol{w} \in W$ such that $w_j \in H_j$ for all $j \in \mathcal{L}$ and identifies $\boldsymbol{u}$ and $\boldsymbol{u}'$ as any two of them.

This procedure enables $D$ to find either both users from $U$ with probability close to 1, or one user from $U$ with probability 1.

*Theorem 4.4:* The family of codes $\mathcal{C} = \{C_k\}$ together with the decision rule given by Algorithm 2 forms a binary fingerprinting code of length $n = mN$ with $2^{Rn} = 2^{K\ell}$ code vectors (users) that either identifies one traitor with probability 1 or both with probability of error

$$p_e \leq 2^{Rn}(2^\ell - 2)^{-N(1-6(1-\delta))}.$$

*Proof:* Let $\boldsymbol{y} = (\boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_N)$ be a fingerprint generated by a coalition $U = \{\boldsymbol{u}, \boldsymbol{u}'\} \subset W$. Suppose that inner decoding of $\boldsymbol{y}$ results in a set $\mathcal{H}$. As we pointed out earlier, for all $i \in \mathcal{I}$ we have $H_i \in \{u_i, u_i'\}$. Hence, for any $\boldsymbol{w} \notin \{\boldsymbol{u}, \boldsymbol{u}'\}$ the number of agreements

$$|\{i : w_i \in H_i\}| \leq |\{i : w_i = u_i\}| + |\{i : w_i = u_i'\}| \leq 2N(1-\delta).$$

Therefore, outer decoding in Step 2 of Algorithm 2 can yield only an element of $U$ (or maybe both of them). If outer decoding gives no output, then $I \leq 4N(1-\delta)$ because

$$I \leq |\{i : H_i = u_i\}| + |\{i : H_i = u_i'\}|.$$

In this case, $L \geq N(1 - 4(1 - \delta))$, and both $u_j, u_j' \in H_j$ for all $j \in \mathcal{L}$.

Now observe that for $\boldsymbol{w} \notin \{\boldsymbol{u}, \boldsymbol{u}'\}$, its coordinate $w_j$ can be included in the list $H_j$ either because $w_j \in \{u_j, u_j'\}$, and the number of such positions is at most $2N(1-\delta)$, or if $w_j \notin \{u_j, u_j'\}$, then it can be included in $H_j$ with probability $p \leq \frac{1}{q-2}$ because the random bijection $\varphi_j^{(k)}$ is not known to the members of $U$, except for the (two) values $\varphi_j^{(k)}(u_j)$ and $\varphi_j^{(k)}(u_j')$. Hence, for all $j \in \mathcal{L}$, the probability

$$\Pr\{w_j \in H_j \mid j \in \mathcal{L}\} \leq (q-2)^{-(L-2N(1-\delta))}$$
$$\leq (q-2)^{-N(1-6(1-\delta))}$$

for any vector $\boldsymbol{w} \in W \backslash \{\boldsymbol{u}, \boldsymbol{u}'\}$. Therefore, the probability that at least one vector from the code $W$ that is distinct from $\boldsymbol{u}, \boldsymbol{u}'$, will satisfy the condition of Step 2 is at most $(2^{K\ell} - 2)(q - 2)^{-N(1-6(1-\delta))}$. This proves the theorem. $\square$

*Example:* Take a $(2, 2)$-separating binary $[35, 6]$ code [21] as inner code $V$, an extended RS $[65, 7]$ code over $\mathbb{F}_{64}$, as outer code $W$, and all $(64!)^{65}$ different vector mappings $\phi^{(k)} = (\varphi_1^{(k)}, \ldots, \phi_{65}^{(k)})$, where the $\varphi_i^{(k)} : \mathbb{F}_{64} \to V$ are bijections. Then Theorem 4.4 states that the resulting binary fingerprinting code of length $n = 2275$ and size $M = 2^{42}$ either identifies one traitor with probability 1 or both with probability of error $p_e < 2^{42} 62^{-29} < 10^{-39}$.

Let us use Theorem 4.4 to construct specific families of codes. Observe that, denoting $\alpha = 1 - 6(1 - \delta)$, we can rewrite the bound for $p_e$ as follows:

$$p_e \leq 2^{-n(R(V)\alpha - R)}(1 - 2^{-\ell+1})^{-N\alpha} < 2^{-n(R(V)\alpha - R)} 4^{2N\alpha/2^\ell} \tag{13}$$

where the last step relies on the inequality $(1 - 1/x)^x > 1/4$ valid for all $x > 2$.

In particular, let us take $W$ an RS code. We obtain the following.

*Corollary 4.5:* For any rate $R < (3 - \log_2 7)/21 \approx 0.0092$, there exists a family of fingerprinting codes $C_k$ of length $n$ with $2^{Rn}$ code vectors (users) that either identify one traitor with probability 1 or both with probability of error $p_e \leq 2^{-n(R_4 - 7R - \beta)}$, where $R_4 = (3 - \log_2 7)/3 \approx 0.0642$ and $\beta = O(N^{-1})$.

*Proof:* Take $V$ a $(2, 2)$ separating $[m, \ell]$ linear code. By Proposition 3.4, for large $m$ there exists a code of rate arbitrarily close to $R_4$. Take $W$, an extended RS $[N = 2^\ell, K]$ code over $\mathbb{F}_q$, $q = N$ of rate $R(W) < 1/7$. The estimate of $p_e$ is obtained from (13) by direct substitution. $\square$

*Remark:* Codes of Corollary 4.5 have a stronger property than codes of the previous section for the case of $t = 2$, namely, the error-free recovery of one user from the coalition or of both users with a small probability of error. This is achieved in exchange for a drop of the maximal achievable rate of codes by a factor of $6/7$ (cf. Corollary 4.5 versus Corollary 4.2).

Now let us concatenate a $(2, 2)$-separating binary $[m, 2\ell]$ code $V$ with a family of AG codes $W$ from maximal curves over $\mathbb{F}_{2^{2\ell}}$. We assume that the rate $R(W)$ is fixed and let $N$ grow. Then the parameters of $W$ approach the bound $1 - \delta = R(W) + 1/(2^\ell - 1)$. The parameters of the resulting code family depend on the choice of the code $V$. For instance, we have the following result.

*Corollary 4.6:* For any rate $R < \frac{121}{8001} \approx 0.015$ there exists a family of fingerprinting codes of length $n$ and rate $R$ that either identify one traitor with probability 1 or both with probability of error $p_e \leq 2^{-n(0.105 - 7R)}$.

*Proof:* Let $V$ be the $(2, 2)$-separating binary $[m = 126, 2\ell = 14]$ code of rate $R(V) = 1/9$ [8] and $W$ be an $[N, K, \delta N]$ linear code from a family of maximal curves over $\mathbb{F}_{2^{14}}$. For large $N$ we have $\alpha = (121/127) - 54R$. Using this in (13) gives the claimed bound on $p_e$. $\qquad\square$

Note that construction complexity of codes given by Corollary 4.6 grows polynomially with the code length $n$.

## V. IDENTIFICATION (DECODING) ALGORITHMS

In this section, we address the algorithmic side of the recovery problem of the users from the coalition. This question was essentially sidestepped in the literature. A notable exception is work on tracing traitors [6], [7] which discusses identification complexity for the narrow-case fingerprinting problem. However, these papers, as well as [5], only give algorithms with complexity linear in the number $M$ of users of the system, i.e., in our terms, exponential in the length of the fingerprinting code. Here, we give algorithms of complexity $\mathsf{polylog}(M) = \mathsf{poly}(n)$ that ensure exponential decline of the error probability.

Conceptually, we face a problem of decoding of two-level concatenated codes. In the case of error correction this problem has a long history in coding theory literature [10], [9]. Decoding can be accomplished, for instance, by exhaustive search performed for the inner-level code and some algebraic algorithm for the outer code. Following this pattern, we analyze the decoding algorithm formed of the following steps.

Recall the parameters $(m, q)$ of the inner code $V$ and $[N, K, \Delta = N\delta]$ of the outer code $W$ which will be assumed an AG code. Here, by $\Delta$ we mean the designed distance of $V$. We will confine ourselves to the so-called one-point AG codes. A one-point code is a geometric Goppa code $W$ constructed in a usual way from a (smooth projective absolutely irreducible) curve over $\mathbb{F}_q$ with rational points $P_0, P_1, \ldots, P_N$. Namely, let $F_1, \ldots, F_K$ be a basis of the Riemann–Roch space of rational functions on $X$ associated with a divisor of the form $aP_0$, $a \geq 0$. Then the set of vectors $\{(F_i(P_1), \ldots, F_i(P_N))$, $1 \leq i \leq K\}$ in $\mathbb{F}_q$ forms a basis of the code $W$. For this reason, all the codes obtainable in this way are sometimes called *evaluation codes*. In particular, if we take $X$ to be the projective line over $\mathbb{F}_q$ and $P_0 = \infty$, then the space of functions is formed by all the polynomials of degree at most $a$. The dimension of this space is $K = a + 1$, and we obtain an $[N, K]$ RS code.

Suppose that a coalition $U = (\boldsymbol{u}^1, \ldots, \boldsymbol{u}^t)$ generates a fingerprint $\boldsymbol{y} = (\boldsymbol{y}_1, \ldots, \boldsymbol{y}_N)$. The distributor first performs decoding of every column $\boldsymbol{y}_i$ with the inner code $V$. The result of

this decoding for the $i$th column is any coalition of size $t$ whose envelope contains $\boldsymbol{y}_i$. Upon inner decoding of all the $\boldsymbol{y}_i$ we obtain a subset

$$\mathcal{H} = \{\boldsymbol{z} \in \mathbb{F}_q^N : z_1 \in H_1, \ldots, z_N \in H_N\}$$

of $N$-words over $\mathbb{F}_q$. The objective of outer decoding is to find a vector $\boldsymbol{w} \in W$ that minimizes the distance $\mathrm{d}(\boldsymbol{w}, \mathcal{H})$ from the code to the "received" subset $\mathcal{H}$. A straightforward approach to this problem requires $|\mathcal{H}|$ runs of some conventional decoding algorithm, which results in exponential running time. Furthermore, until recently all known algebraic algorithms could correct only about $\Delta/2$ errors (the so-called bounded distance decoding). Note that even if the set $\mathcal{H}$ is of size 1, it does not help to use a bounded distance decoding algorithm because we need to correct, roughly speaking, $N(1 - 1/t) \geq N/2$ errors. To show this, recall that users in the coalition can create any vector $\boldsymbol{y}$ that is contained in their envelope. Consider the following fingerprint vector $\boldsymbol{y} = (\boldsymbol{y}_1, \ldots, \boldsymbol{y}_N)$, where (assuming that $N/t$ is integer)

$$\boldsymbol{y}_{st+i} = \varphi_{st+i}^{(k)}(u_{st+i}^i), \quad s = 0, 1, \ldots, N/t - 1, \ i = 1, \ldots, t.$$

Then, $H_{st+i} = u_{st+i}^i$, so $|H_{st+i}| = 1$, and for every vector $\boldsymbol{u}^i$ we have $N - \mathrm{d}(\boldsymbol{u}^i, \mathcal{H}) = N/t$. Fortunately, the recently found list decoding procedures of AG codes correct many more errors, outputting in polynomial time a small (polynomial-sized) list of code vectors. More precisely, we have Guruswami–Sudan (GS) decoding [11], [12]. Let $W$ be a $q$-ary (one-point) AG code with parameters $[N, RN, \delta N]$. Let $\boldsymbol{R} = [r_{ji}]$ be a $q \times N$ matrix of nonnegative integers. There exists an algorithm which finds all the codewords $\boldsymbol{w} \in W$ that satisfy the inequality

$$r(\boldsymbol{w}) := \sum_{i=1}^{N} r_{w_i, i} \geq \sqrt{N - \Delta} \, \|\boldsymbol{R}\| \qquad (14)$$

where $\|\boldsymbol{R}\| = \sqrt{\sum_{i, j} r_{ji}^2}$ is the $\ell_2$-norm of $\boldsymbol{R}$. The number of these codewords is bounded by a polynomial function of $N$. The implementation complexity of the algorithm is also polynomial in $N$. Implementation details for the RS and one-point AG codes together with complexity estimates are worked out in [17], [12].

This formulation of the decoding algorithm is motivated by the soft-decision decoding setting. In that context, the matrix $\boldsymbol{R}$ describes reliabilities of symbols of $\mathbb{F}_q$. For this reason, in the following we call the value $r(\boldsymbol{w})$ the reliability of the vector $\boldsymbol{w}$. The original list decoding algorithm of [11] that works with a particular "received" vector $\boldsymbol{x} = (x_1, x_2, \ldots, x_N)$ is obtained by putting $r_{ji} = \delta_{j, x_i}$ (ones for the values that correspond to the $x_i$ and zeros elsewhere). In this case, by (14), the algorithm will find a list $\mathcal{Y} \subset W$ formed of all the codewords $\boldsymbol{w}$ of the code $W$ that satisfy $\mathrm{d}(\boldsymbol{w}, \boldsymbol{x}) \leq N - \sqrt{N(N - \Delta)}$. The size $|\mathcal{Y}|$ is again bounded by a polynomial function of $N$.

Note that if our goal is to output one decoding result rather than a list, we choose a vector $\boldsymbol{w}$ from $\mathcal{Y}$ whose reliability $r(\boldsymbol{w})$ is the maximum of all the code vectors in $\mathcal{Y}$. In particular, for the basic version of the algorithm which builds a list of code vectors in a sphere of radius $N - \sqrt{N(N - \Delta)}$ around a given point $\boldsymbol{x}$, we can choose a vector whose distance to $\boldsymbol{x}$ is the smallest among all the members of $\mathcal{Y}$. In other words, this algorithm

enables one to correct, whenever possible, a number of errors greater than half the minimum distance of the code.

This observation explains the usefulness of the list decoding methods in the fingerprinting problem. In the following subsections, we will show a way to find a member of the coalition as a most reliable code vector, maintaining the overall polynomial complexity.

### A. Identification in the Case of Arbitrary $t$

Consider the family of binary fingerprinting codes $\mathcal{C}$ formed by concatenating a binary $(t, t)$-separating $(m, q)$ code $V$ of rate $R(V) = m^{-1} \log_2 q$ with outer $q$-ary $[N, R(W)N, \delta N]$ AG codes $W$.

Consider the following identification procedure for the observed fingerprint $\boldsymbol{y} = (\boldsymbol{y}_1, \ldots, \boldsymbol{y}_N)$.

Algorithm 1′:

1) For every $i = 1, 2, \ldots, n$, find the subset $H_i$ of symbols of $\mathbb{F}_q$ corresponding to an arbitrary $t$-tuple of vectors of $V$ such that $\boldsymbol{y}_i$ is contained in their envelope. To accomplish this, the distributor performs a lookup of at most $\binom{m}{t}$ $t$-tuples of vectors of $V$. (Same as Step 1 of Algorithm 1.)

2) Form a $(q \times N)$ matrix $\boldsymbol{R}$, setting for $i = 1, 2, \ldots, N$
$$r_{ji} = 1, \qquad \text{if } j \in H_i$$
$$r_{ji} = 0, \qquad \text{if } j \notin H_i.$$

3) Use the GS algorithm with the code $W$ and matrix $\boldsymbol{R}$ as the input to obtain a list $\mathcal{Y} \subset W$ of code vectors of $W$ that satisfy (14). If $\mathcal{Y} \neq \emptyset$, find a vector $\boldsymbol{u}$ from $\mathcal{Y}$ such that the value $r(\boldsymbol{u}) \geq N/t$. Identify $\boldsymbol{u}$ as a member of the coalition $U$.

*Remark:* The reliabilities $r_{ji}$ are chosen in a way suitable for addressing the following decoding problem: given a subset $\mathcal{H} \subset \mathbb{F}_q^N$
$$\mathcal{H} = (\boldsymbol{h} \in \mathbb{F}_q^N : h_1 \in H_1, \ldots, h_N \in H_N)$$
find a codevector of the code $W$ that minimizes the Hamming distance between $W$ and $\mathcal{H}$. This generalizes the standard decoding problem for which $|\mathcal{H}| = 1$.

*Theorem 5.1:* Let $\mathcal{C} = \{C_k\}$ be a family of binary fingerprinting codes of rate $R = R(W)R(V)$ with $(t, t)$-separating inner code of rate $R(V) = m^{-1} \log_2 q$ and outer $q$-ary code $W$ of rate $R(W)$ and relative distance $\delta > 1 - 1/t^3$.

For any $\gamma > 0$, there exists a value $q_0 = q_0(\gamma)$ such that for any $q \geq q_0$ and any $t \leq (1 - \gamma)\sqrt{q - 1}$, Algorithm 1′ used in conjunction with the codes $C_k$ has complexity $\text{poly}(n)$ and identifies a member of the coalition with error probability
$$p_e \leq 2^{-N[D(\sigma \| p) - t^{-3} \log_2 q]}$$
where $\sigma = (t - 1)/t^2$, $p = (t - 1)/(q - 1)$.

*Proof:* By definition, for every $i$ the subset $H_i$ contains at least one of the coordinates of the vectors of the coalition $U$. Since for any $\boldsymbol{u}$, the value $r(\boldsymbol{u}) = |\{i : u_i \in H_i\}|$, we have
$$\sum_{\boldsymbol{u} \in U} r(\boldsymbol{u}) \geq N$$
and, therefore, there exists a vector $\hat{\boldsymbol{u}} \in U$ such that $r(\hat{\boldsymbol{u}}) \geq N/t$. We shall prove that

i) $\hat{\boldsymbol{u}}$ satisfies (14), and

ii) the probability $p_e$ of incorrect identification satisfies the claim of the theorem.

We begin with part i). By assumption, $r(\hat{\boldsymbol{u}}) \geq N/t$. On the other hand, the right-hand side of (14) for our choice of the matrix $\boldsymbol{R}$ equals
$$\sqrt{(N - \Delta) \sum_{i=1}^{N} \sum_{j=1}^{q} r_{ji}^2} = \sqrt{(N - \Delta)tN}$$
$$= N\sqrt{(1 - \delta)t} \leq \frac{N}{t}.$$

Hence, the vector $\hat{\boldsymbol{u}} \in U$ satisfies (14).

For part ii), note that
$$p_e \leq \Pr\{\exists \boldsymbol{w} \in W \backslash U : r(\boldsymbol{w}) \geq N/t\}.$$

Therefore, let us estimate the probability that there exists a vector $\boldsymbol{w} \in W \backslash U$ such that $r(\boldsymbol{w}) \geq N/t$. Denote by $I$ the number of coordinates $i$ in which $w_i$ agrees with one of the vectors in $U$. Assume that these coordinates have numbers $1, 2, \ldots, I$. We have
$$I \leq t(N - \Delta) \leq \frac{N}{t^2}.$$

If $j$ is one of the remaining $N - I$ coordinates, then by (9) we have $\Pr(w_j \in H_j) \leq p = \frac{t-1}{q-1}$. Therefore, the probability of mistakenly identifying a vector $\boldsymbol{w}$ is estimated as follows:
$$P_{e, \boldsymbol{w}} = \Pr\left\{ r(\boldsymbol{w}) \geq \frac{N}{t} \right\} \leq \Pr\left\{ I + \sum_{j=I+1}^{N} r_{w_j, j} \geq \frac{N}{t} \right\}$$
$$\leq \Pr\left\{ \sum_{i=1}^{N} \xi_i \geq N \frac{t-1}{t^2} \right\}$$

where every $\xi_i$ is a Bernoulli random variable that takes the value 1 with probability $p$ and 0 with probability $1 - p$. The overall error probability of identification can be bounded as $p_e \leq q^{NR(W)} P_{e, \boldsymbol{w}}$. This quantity can be estimated as follows:
$$p_e \leq 2^{-N(D(\sigma \| p) - R(W) \log_2 q)} \leq 2^{-N(D(\sigma \| p) - t^{-3} \log_2 q)}.$$

To complete the proof of the theorem, we have to show that for our choice of the parameters the exponent
$$E = D(\sigma \| p) - t^{-3} \log_2 q$$
of this bound is a positive number. The exponent has the form
$$E(q, t) = \log_2 \frac{q-1}{t^2} + \frac{t^2 - t + 1}{t} \log_2 \frac{t^2 - t + 1}{q - t} - \frac{1}{t^3} \log_2 q.$$

Its derivative $\partial E / \partial t$ can be checked to be negative for all $t \leq \sqrt{q - 1}$. Therefore, it suffices to check that for any $\gamma > 0$ there exists a value $q_0$ such that $E[q, (1 - \gamma)\sqrt{q - 1}] > 0$ for any $q \geq q_0$. The quantity
$$\frac{1}{t^3} \log_2 q = \frac{\log_2 q}{((q - 1)^{1/2}(1 - \gamma))^3}$$

can be made arbitrarily small for large $q$. The term $D(\sigma \| p)$ is zero for $t^2 = q - 1$ and positive otherwise. For any fixed $t$ this

term is a growing function of $q$. Hence, for any $\gamma > 0$ there exists a value $q_0 = q_0(\gamma)$ such that for any $q > q_0$

$$E[q, (1 - \gamma)\sqrt{q - 1}] > 0.$$

This completes the proof. $\qquad\square$

The value $q_0(\gamma)$ is easy to find numerically. For instance, $q_0(1/3) = 39$, $q_0(1/10) = 355$, $q_0(1/20) = 1613$, etc.

This theorem is valid for any sequence of AG codes (including the RS codes). To maximize the achievable rate of the codes $\mathcal{C}$ we can choose a sequence of AG codes [27] whose parameters for large $N$ approach the bound

$$R(W) = 1 - \delta - \frac{1}{\sqrt{q} - 1}.$$

This shows that for a given $(t, t)$ separating code $V$ of length $m$ and size $q$ there exist sequences of $t$-secure fingerprinting codes $\{C_k\}$ with rate $R = R(V)R(W)$, where

$$0 \leq R(V) \leq -\frac{1}{2t - 1}\log_2(1 - 2^{2t-1}) - \frac{1}{m}$$

$$0 \leq R(W) \leq \frac{1}{t^3} - \frac{1}{\sqrt{q} - 1}.$$

The error probability of identification for codes $C_k$ falls exponentially in the code length $n$. Hence, taking $q$ sufficiently large we derive the following result.

*Corollary 5.2:* Let $R_t^{(s)}$ be the maximum achievable rate of $(t, t)$-separating codes. For any rate $R$, $0 < R < R_t^{(s)}/t^3$, there exists a sequence of $t$-secure fingerprinting codes of length $n$ and size $2^{Rn}$ that allow polynomial-time identification with error probability falling exponentially with the code length $n$.

*B. Identification Algorithm in the Case $t = 2$*

The case $t = 2$ is of special interest because it is possible to modify the general decoding algorithm of the last section so that it lends itself to a more accurate performance analysis.

Consider the family of fingerprinting codes $\mathcal{C}$ of length $n = Nm$ and size $2^{\ell N R(W)}$ formed by concatenating a binary $(2, 2)$-separating $[m, \ell]$ code $V$ ($\ell$ is even) and a $q$-ary $[N, R(W)N, \delta N]$ AG code $W$ ($q = 2^\ell$). By Section IV-B, after inner decoding the distributor forms a set

$$\mathcal{H} = \{\boldsymbol{h} \in \mathbb{F}_q^N : h_1 \in H_1, \ldots, h_N \in H_N\}$$

where some subsets $H_i$ are singletons (i.e., $q$-ary letters) and some 3-sets of $q$-ary letters. We recall the notation $\mathcal{I}$ for the subset of coordinates corresponding to the singletons and $\mathcal{L} = \{1, 2, \ldots, N\}\backslash\mathcal{I}$ for the subset corresponding to the triplets; $I = |\mathcal{I}|$, $L = |\mathcal{L}|$.

Consider the following decoding procedure of the code $C = C_k$.

Algorithm 2':
1) Apply the same procedure as in Step 1 of Algorithm 2 to construct the set $\mathcal{H}$.
2) Form the matrix $\boldsymbol{R}$ as follows:

$$r_{\alpha, i} = \begin{cases} 2, & \alpha = h_i \\ 0, & \alpha \neq h_i \end{cases} \qquad r_{\alpha, j} = \begin{cases} 1, & \alpha \in H_j \\ 0, & \alpha \notin H_j \end{cases}$$

where $i \in \mathcal{I}$ and $j \in \mathcal{L}$.
3) Use the GS algorithm with the code $W$ and matrix $\boldsymbol{R}$ as the input to obtain a list $\mathcal{Y} \subset W$ of code vectors $\boldsymbol{w}$ that satisfy

$$r(\boldsymbol{w}) \geq N.$$

If $\mathcal{Y} \neq \emptyset$, identify an arbitrary vector from $\mathcal{Y}$ as a member of the coalition.

The difference of this algorithm from the general Algorithm 1' is motivated by the fact that the set $\mathcal{H}$ for the case $t = 2$ has a particular structure described above.

*Theorem 5.3:* Let $\mathcal{C} = \{C_k\}$ be a family of binary fingerprinting codes of rate $R = R(W)R(V)$ with $(2, 2)$-separating inner code of rate $R(V) = m^{-1}\log_2 q$ and outer $q$-ary code $W$ of length $N$, rate $R(W)$, and relative distance $\delta \geq 3/4$.

For the family of fingerprinting codes $\mathcal{C}$, the identification rule given by Algorithm 2' has complexity $\mathsf{poly}(n)$ and identifies a member of the coalition with probability $1 - p_e$, where

$$p_e \leq \left(\frac{q - 2}{4\sqrt{q}}\right)^{-N/2}.$$

*Proof:* Let $\boldsymbol{u}$, $\boldsymbol{u}'$ be the vectors of $W$ that correspond to the members of the coalition $U$. At least one of these vectors, say $\boldsymbol{u}$, satisfies the condition $|\{i \in \mathcal{I}: u_i = h_i\}| \geq I/2$. Therefore,

$$\sum_{i=1}^{N} r_{u_i, i} \geq L + 2(I/2) = N.$$

On the other hand

$$(N - \Delta)\|\boldsymbol{R}\|^2 \leq \frac{1}{4}N(4I + 3L) = \frac{1}{4}N(3N + I) \leq N^2.$$

Hence, by (14), the vector $\boldsymbol{u}$ will be contained in the list $\mathcal{Y}$.

Now let $\boldsymbol{w} \in W\backslash U$ be another code vector. The number of coordinates $i$ such that $w_i \in \{u_i, u_i'\}$ is not more than $2(N - \Delta)$. In these coordinates, the value $r_{w_i, i} = 2$ or 1 according as $i \in \mathcal{I}$ or $i \in \mathcal{L}$. In the remaining coordinates, which necessarily fall in the subset $\mathcal{L}$, the value $r_{w_i, i} = 1$ with probability $1/(q - 2)$. So denoting by $\xi_i$ a Bernoulli random variable that takes the value 1 with probability $1/(q-2)$ and 0 with probability $1 - 1/(q - 2)$, we can write the reliability of the vector $\boldsymbol{w}$ as follows:

$$r(\boldsymbol{w}) \leq 2(N - \Delta) + \sum_{i=1}^{N-2(N-\Delta)} \xi_i$$

$$\leq 2(N - \Delta) + \sum_{i=1}^{N} \xi_i \leq \frac{N}{2} + \sum_{i=1}^{N} \xi_i.$$

The probability $P_{e, \boldsymbol{w}}$ of erroneously identifying $\boldsymbol{w}$ can be bounded as follows:

$$P_{e, \boldsymbol{w}} \leq \Pr\{\boldsymbol{w} \in \mathcal{Y}\} = \Pr\{r(\boldsymbol{w}) \geq N\} \leq \Pr\left\{\sum_{i=1}^{N} \xi_i \geq \frac{N}{2}\right\}$$

$$\leq 2^{-ND(\frac{1}{2}\|\frac{1}{q-2})} = 2^{\frac{1}{2}N\log_2\frac{4(q-3)}{(q-2)^2}}$$

$$\leq 2^{\frac{1}{2}N\log_2\frac{4}{q-2}} = \left(\frac{q - 2}{4}\right)^{-N/2}.$$

The size of the code $W$ is $q^{R(W)N} < q^{N(1-\delta)} \le q^{N/4}$. Therefore, we finally obtain for the error probability of identification the bound

$$p_e \le |W| p_{e,\boldsymbol{w}} < q^{N/4} \left(\frac{q-2}{4}\right)^{-N/2} = \left(\frac{q-2}{4\sqrt{q}}\right)^{-N/2},$$

as was to be proved. $\qquad\square$

*Remark:* Codes secure against coalitions of size two were proven in Section IV-B to be able to restore either one user with probability one or both with an exponentially small error probability. With a little more work, it is possible to attain this performance with polynomial complexity. The details will be omitted.

The remarks made after the proof of Theorem 5.1 are also valid for the result proved in this section. In particular, let us combine the $(2, 2)$-separating binary $[m = 126, \ell = 14]$ code with a sequence of AG codes over the field $\mathbb{F}_{2^{14}}$ constructed from maximal curves. We obtain the following result.

*Corollary 5.4:* For any rate $R \le 0.0242$ there exists a sequence of 2-secure fingerprinting codes of length $n$ and size $2^{Rn}$ that allow polynomial-time identification of one of the members of the coalition with error probability $p_e < 2^{-5n/252}$.

This result follows immediately from Theorem 5.3. In particular, since $n = mN$, the error probability is at most

$$p_e < \left(\frac{q-2}{4\sqrt{q}}\right)^{-N/2} \approx 2^{-\frac{5}{2}N} = 2^{-5n/252}.$$

It is interesting to compare this result with Corollary 4.6. We see that the exponent of the error probability in that corollary is better than here for rates $R$ close to zero. On the other hand, Corollary 5.4 gives a constructive family of codes with rates much greater than in Corollary 4.5. Moreover, the error probability of the codes just constructed is uniformly bounded for all the values of the rate $R$, whereas in Corollary 4.5 the error exponent depends on $R$ and approaches zero when the rate $R$ approaches the maximum admissible value.

## VI. CONCLUDING REMARKS: RELATED PROBLEMS

We considered properties and constructions of fingerprinting codes, mainly for the binary case. Here we wish to discuss two problems studied in the literature (see [24]) that are related to our paper, and to state an open problem.

1) *Codes with the identifiable parent property (t-IPP codes) and t-traceability codes.* As remarked in Section II-A, codes for the narrow-sense fingerprinting problem with zero error are also called IPP codes. This problem was suggested in [13], where it was also proved that there exist $q$-ary 2-IPP codes with positive rate for every $q \ge 3$. The question of code existence for other $t$ was discussed in [24] and settled in full in [1]: there exist $t$-IPP codes with rate $R > 0$ if and only if $q \ge t + 1$.

A $t$-IPP code $C$ has the traceability property [6], [7] if for every fingerprint $\boldsymbol{y}$ created by the coalition $U$, the set of its nearest neighbors in $C$ is contained in $U$. It is known [6] that a $q$-ary $[n, k, d]$ code possesses the $t$-traceability property if $d \ge n(1 - 1/t^2)$. It is clear from the properties of the GS algorithm that evaluation (AG or RS) codes with large distance pos-

sess the traceability property together with a polynomial identification algorithm. (After this paper was submitted, we became aware of the paper [23] which works out the details of this idea.)

This application of list decoding is, however, limited as follows: to construct $q$-ary error-correcting codes of rate bounded away from zero, we must take $q > n/(n - d)$ (by the Plotkin bound of coding theory). Hence, if traceability codes or efficient IPP codes are constructed based on error-correction properties only, then no matter what the identification algorithm is, nonzero code rate is obtained only if $q > t^2$.

On the other hand, the ideas of the present paper, in a simplified form, enable one to construct explicit families of $t$-IPP codes together with polynomial-time identification procedures for any $q \ge t + 1$ (this was suggested as an open problem in [24]). This result is presented elsewhere [2], [3]. Independently, similar results were obtained in [26].

2) *Capacity of the fingerprinting channel* (open problem). Let us call the rate $R$ an achievable rate of the "fingerprinting channel" if for any given $a > 0$ there exists a positive integer $n_0$ with the property that for every $n > n_0$ the best attainable error probability of identification (2) satisfies

$$\varepsilon(\mathcal{E}; n, q^{nR}, t) < a.$$

As usual, let us call the number

$$\mathcal{C}(\mathcal{E}; q, t) = \sup_{R \text{ achievable}} R$$

capacity of the digital fingerprinting channel for the envelope $\mathcal{E}$. We have shown that $\mathcal{C}(\mathcal{E}; 2, t) > 0$ for every given $t$ and presented constructive code sequences with rate bounded away from zero. The question of finding the capacity of the fingerprinting channel presents an interesting open problem.

### REFERENCES

[1] A. Barg, G. Cohen, S. Encheva, G. Kabatiansky, and G. Zémor, "A hypergraph approach to the identifying parent property: The case of multiple parents," *SIAM J. Discr. Math.*, vol. 14, pp. 423–431, 2001.

[2] A. Barg and G. Kabatiansky, "A class of IPP codes with efficient identification," in *Proc. 39th Annual Allerton Conf. Communications, Control and Computing*, Monticello, IL, Oct. 2001, pp. 885–890.

[3] ——, "A class of IPP codes with efficient identification," DIMACS Tech. Rep. 2002-36, submitted for publication.

[4] G. R. Blakley, C. Meadows, and G. Purdy, "Fingerprinting long forgiving messages," in *Proc. CRYPTO*, 1985, pp. 180–189.

[5] D. Boneh and J. Shaw, "Collusion-secure fingerprinting for digital data," *IEEE Trans. Inform. Theory*, vol. 44, pp. 1987–1905, Sept. 1998.

[6] B. Chor, A. Fiat, and M. Naor, "Tracing traitors," in *Crypto'94 (Lecture Notes in Computer Science)*. Berlin, Germany: Springer-Verlag, 1994, vol. 839, pp. 257–270.

[7] B. Chor, A. Fiat, M. Naor, and B. Pinkas, "Tracing traitors," *IEEE Trans. Inform. Theory*, vol. 46, pp. 893–910, May 2000.

[8] G. Cohen, S. Encheva, and S. Litsyn, "Intersecting codes and partially identifying codes," in *Proc. Int. Workshop on Coding and Cryptography*, Ecoles de Coëtquidan, Paris, France, 2001, pp. 139–147.

[9] I. Dumer, "Concatenated codes and their multilevel generalizations," in *Handbook of Coding Theory*, V. Pless and W. C. Huffman, Eds. Amsterdam, The Netherlands: Elsevier Science, 1998, vol. 2, pp. 1911–1988.

[10] G. D. Forney, Jr., *Concatenated Codes*. Cambridge, MA: MIT Press, 1966.

[11] V. Guruswami and M. Sudan, "Improved decoding of Reed–Solomon and algebraic-geometry codes," *IEEE Trans. Inform. Theory*, vol. 45, pp. 1757–1767, Sept. 1999.

[12] V. Guruswami, "List decoding of error-correcting codes," Ph.D. dissertation, MIT, Cambridge, MA, 2001.

[13] H. D. Hollmann, J. H. van Lint, J.-P. Linnartz, and L. M. G. M. Tol-huizen, "On codes with the identifiable parent property," *J. Combin. Theory*, ser. A, vol. 82, pp. 121–133, 1998.

[14] A. D. Friedman, R. L. Graham, and J. D. Ullman, "Universal single transition time asynchronous state assignments," *IEEE Trans. Comput.*, vol. C–18, pp. 541–547, 1969.

[15] J. Körner and G. Simonyi, "Separating partition systems and locally different sequences," *SIAM J. Discr. Math.*, pp. 355–359, 1988.

[16] F. J. MacWilliams and N. J. A. Sloane, *The Theory of Error-Correcting Codes*. Amsterdam, The Netherlands: North-Holland, 1977.

[17] R. R. Nielsen, "Decoding AG-codes beyond half the minimum distance," Master thesis, Dept. Math., Tech. Univ. Denmark, Lyngby, 1998.

[18] M. S. Pinsker and Y. L. Sagalovich, "Lower bound on the cardinality of code of automata's states," *Probl. Inform. Transm.*, vol. 8, no. 3, pp. 59–66, 1972.

[19] Y. L. Sagalovich, "Concatenated codes of the states of an automaton," *Probl. Inform. Transm.*, vol. 14, no. 2, pp. 77–85, 1978.

[20] ——, "Fully separated systems," *Probl. Inform. Transm.*, vol. 18, no. 2, pp. 74–82, 1982.

[21] ——, "Separating systems," *Probl. Inform. Transm.*, vol. 30, no. 2, pp. 14–35, 1994.

[22] K. W. Shum, I. Aleshnikov, P. V. Kumar, H. Stichtenoth, and V. Deolalikar, "A low-complexity algorithm for the construction of algebraic-geometric codes better than the Gilbert–Varshamov bound," *IEEE Trans. Inform. Theory*, vol. 47, pp. 2225–2241, Sept. 2001.

[23] A. Silverberg, J. Staddon, and J. L. Walker, "Efficient traitor tracing algorithm using list decoding," in *ASIACRYPT 2001 (Lecture Notes in Computer Science)*. Berlin, Germany: Springer-Verlag, 2001, vol. 2248, pp. 175–192.

[24] J. N. Staddon, D. R. Stinson, and R. Wei, "Combinatorial properties of frameproof and traceability codes," *IEEE Trans. Inform. Theory*, vol. 47, pp. 1042–1049, Mar. 2001.

[25] D. R. Stinson, V. T. Tran, and R. Wei, "Secure frameproof codes, key distribution patterns, group testing algorithms and related structures," *J. Stat. Plan. Inference*, vol. 86, pp. 595–617, 1998.

[26] T. V. Tran and S. Martirosian, "Constructions for efficient IPP codes," preprint, 2002.

[27] M. Tsfasman and S. Vlăduţ, *Algebraic-Geometric Codes*. Dordrecht, The Netherlands: Kluwer, 1991.