

A Capacity-Achieving Simple Decoder for Bias-Based Traitor Tracing Schemes

Jan-Jaap Oosterwijk, Boris Škorić, Jeroen Doumen

Abstract—We investigate alternative suspicion functions for bias-based traitor tracing schemes, and present a practical construction of a simple decoder that attains capacity in the limit of large coalition size c .

We derive optimal suspicion functions in both the Restricted-Digit Model and the Combined-Digit Model. These functions depend on information that is usually not available to the tracer – the attack strategy or the tallies of the symbols received by the colluders. We discuss how such results can be used in realistic contexts.

We study several combinations of coalition attack strategy versus suspicion function optimized against some attack (another attack or the same). In many of these combinations the usual code length scaling $\ell \propto c^2$ changes to a lower power of c , e.g. $c^{3/2}$. We find that the interleaving strategy is an especially powerful attack. The suspicion function tailored against interleaving is the key ingredient of the capacity-achieving construction.

Index Terms—Traitor tracing, collusion resistance.

I. INTRODUCTION

A. Collusion attacks on watermarking

FORENSIC watermarking is a means for tracing the origin and distribution of digital content. Before distribution, the content is modified by embedding an imperceptible watermark, which plays the role of a personalized serial number. Once an unauthorized copy of the content is found, the identities of those users who participated in its creation can be determined. A tracing algorithm outputs a list of suspicious users.

The most powerful attacks against watermarking are *collusion attacks*, in which multiple attackers (the ‘coalition’) combine their differently watermarked versions of the same content; the observed differences point to the locations of the hidden marks.

In the past two decades several types of collusion-resistant codes have been developed. The most popular type in the recent literature is the class of *bias-based* codes. These were introduced by G. Tardos in 2003. The original paper [1] was followed by a flurry of activity, e.g. improved analyses [2]–[7], code modifications [8]–[10], decoder modifications [11]–[13] and various generalizations [14]–[17]. The advantage of bias-based versus deterministic codes is that they can achieve the asymptotically optimal relationship $\ell \propto c^2$ between the sufficient code length ℓ and the coalition size c .

Two kinds of tracing algorithm can be distinguished: (i) *simple decoders*, which assign a level of suspicion to single users

This research is supported by the Dutch Technology Foundation STW, which is part of the Netherlands Organisation for Scientific Research (NWO), and which is partly funded by the Ministry of Economic Affairs.

B. Škorić are with the Eindhoven University of Technology; J. Doumen is with Irdeto B.V.; J. Oosterwijk is with both.

and (ii) *joint decoders* [11]–[13], which look at sets of users. Joint decoders employ a simple decoder as a bootstrapping step.

Tardos’ scheme worked with a binary code and a simple decoder. Its ‘suspicion function’ for computing a level of suspicion for single users was improved [15] and the scheme was generalized to q -ary alphabets. However, it turns out [18] that the suspicion function yields sub-optimal fingerprinting rates for $q > 3$, i.e. rather far below the fingerprinting capacity [19], [20] and far below the dynamic code rate [21].

Alternative suspicion functions for the binary case were introduced [12], where an Expectation Maximization (EM) algorithm was used. A candidate coalition is selected, which (if the guess is sufficiently good) makes it possible to estimate the employed attack strategy; a suspicion function is then used which is optimized against that strategy. This leads to a new ranking of users, giving a new candidate coalition, and the whole process is repeated until it converges.

B. Contributions

This paper is an extended version of earlier work on optimal suspicion functions [22].

- We generalize the work of Charpentier et al. [12] to q -ary alphabets. Using functional derivation methods we obtain suspicion functions that for large c maximize the expected score for the coalition, allowing the tracer to distinguish best between them and the innocent users. We present results for the Combined-Digit Model and the Restricted-Digit Model.
- We consider a set of often-considered attack strategies. We substitute these attacks into the generic formulas and obtain closed-form expressions for the optimal suspicion functions associated with these attacks.
- We tabulate the performance for each combination of attack and suspicion function. For some cases we prove theorems analytically and for all binary cases we have numerical results. Naturally, in case of a match the sufficient code length ℓ is small; for all considered strategies but the interleaving attack we even find $\ell \propto c^{3/2}$. For the interleaving attack and its matching suspicion function we find an asymptotic fingerprinting rate $(q - 1)/(2c^2 \ln q)$, which is exactly the q -ary asymptotic fingerprinting capacity.

In non-matching cases the results differ widely. In some cases, as expected, the mismatched defense fails completely, while in others the code length remains $\ell \propto c^2$ (often with a smaller coefficient than with the Tardos

suspicion function), and in many cases we find $\ell \propto c^{3/2}$ even for a mismatch.

- The suspicion function tailored against the interleaving attack is very special. When this suspicion function is adopted as the basis of a simple decoder, the minimax game for the asymptotic code rate (attack strategy versus bias distribution function) has a saddle point when the interleaving attack is used and the distribution function is the Dirichlet distribution with concentration parameter $1/2$. In the saddle point the asymptotic rate equals the asymptotic capacity. The saddle point is the same point that was found by Huang and Moulin[20] for the mutual information minimax game. Thus, we have identified a simple decoder that asymptotically achieves capacity.

In Sections III-A and XI we comment on possible ways to exploit our results for the construction of improved decoders by using several suspicion functions in parallel, and/or deploying a tally-dependent suspicion to strengthen the EM algorithm, and/or to validate candidate coalitions in general.

II. PRELIMINARIES

A. General notation

We denote random variables by capital letters and their realizations in lower case. We write vectors in boldface. We define $[\ell] = \{1, \dots, \ell\}$. The q -ary alphabet is \mathcal{A} , which is sometimes set to $\mathcal{A} = \{0, \dots, q-1\}$.

We use multi-index notation, e.g. $\mathbf{p}^\kappa = \prod_{\alpha \in \mathcal{A}} p_\alpha^\kappa$, $\mathbf{p}^m = \prod_{\alpha \in \mathcal{A}} p_\alpha^{m_\alpha}$, and $\binom{c}{\mathbf{m}} = c! / \prod_{\alpha \in \mathcal{A}} m_\alpha!$.

We define the norm of a vector as $|\mathbf{p}| = \sum_{\alpha \in \mathcal{A}} p_\alpha$. For probability mass/density functions we use abbreviated notation of the form $f_{y|\mathbf{p}} = f_{Y|\mathbf{P}}(y|\mathbf{p})$ when it does not cause ambiguity.

In conditional expectation values we sometimes use the abbreviation $\mathbb{E}_{\mathbf{M}|\mathbf{p}}[\dots] = \mathbb{E}_{\mathbf{M}}[\dots | \mathbf{P} = \mathbf{p}]$. An \mathbb{E} without subscripts is an expectation over *all* probabilistic degrees of freedom. We use $\delta_{x,y}$ to denote the Kronecker delta function, which is 1 when $x = y$ and 0 when $x \neq y$.

The notation $\frac{\partial A}{\partial p_x}|_{|\mathbf{p}|=1}$ is defined as follows. First the derivative $\partial A / \partial p_x$ is taken *without* taking the constraint $\sum_{\alpha} p_\alpha = 1$ into account. After differentiation the constraint is enforced.

We will use the shorthand notation $a_k := (p_0 + \dots + p_{k-1})$ and $a_B = \sum_{\beta \in B} p_\beta$.

B. Bias-based tracing; simple decoder

The content contains ℓ abstract ‘locations’ into which a q -ary symbol can be embedded. For each location $i \in [\ell]$ independently, the tracer draws a bias vector $\mathbf{P}_i = (P_{i,\alpha})_{\alpha \in \mathcal{A}}$ from a distribution $f_{\mathbf{P}}$. The biases satisfy $P_{i,\alpha} \geq 0$ and $|\mathbf{P}_i| = 1$. A symmetric Dirichlet distribution was taken [15], with concentration parameter $\kappa > 0$,

$$f_{\mathbf{P}}(\mathbf{p}) = \mathbf{p}^{\kappa-1} \Gamma(q\kappa) / [\Gamma(\kappa)]^q. \quad (1)$$

For $q = 2$ it is customary to set $\kappa = \frac{1}{2}$, turning (1) into the arcsine distribution for the component p_1 . However, in that case the support has to be reduced to $p_1 \in [\delta, 1-\delta]$, with cutoff

parameter $\delta > 0$, in order to avoid statistical problems due to extremely unlikely events. The probability density function then becomes

$$f_{\mathbf{P}}(p_1) = \frac{1}{2 \arcsin(1-2\delta)} \frac{1}{\sqrt{p_1(1-p_1)}}. \quad (2)$$

As the cutoff parameter is typically chosen so small that it vanishes, we will neglect it in our analysis. The number of users is n . For each $i \in [\ell]$ and each $j \in [n]$, the tracer draws a random symbol $X_{i,j} \in \mathcal{A}$ according to the categorical distribution \mathbf{P}_i , i.e. $\mathbb{P}[X_{i,j} = \alpha | \mathbf{P}_i = \mathbf{p}_i] = p_{i,\alpha}$ independent of j . The symbol $X_{i,j}$ is embedded into the content of user j in location i .

The coalition of attackers is denoted as $\mathcal{C} \subset [n]$, with $|\mathcal{C}| = c$. In some attack models, e.g. the Combined-Digit Model (Section II-C), they are allowed to do signal processing attacks such as introducing noise and fusing symbols. In the Restricted-Digit Model (RDM) they are only allowed to select one colluder’s symbol (denoted as y_i) in location i . In the *simple decoder* approach, the tracer determines a score S_j for each user j by adding independently computed sub-scores $S_{i,j}$ for each location i ; these are based on \mathbf{p}_i , $X_{i,j}$ and the colluders’ output in location i . If the score exceeds a threshold, user j is accused.

Tardos [1] introduced a (simple decoder) score system for the RDM at $q = 2$ that was later [15] symmetrized and generalized to $q > 2$. The sub-scores for each location are computed using a ‘suspicion function’ g as $S_{i,j} = g(x_{i,j}, y_i, \mathbf{p}_i)$ with

$$g(x, y, \mathbf{p}) = \begin{cases} \sqrt{(1-p_y)/p_y} & \text{if } x = y \\ -\sqrt{p_y/(1-p_y)} & \text{if } x \neq y. \end{cases} \quad (3)$$

It has the special property that the $S_{i,j}$ of innocent users has expectation 0 and variance 1.

Given the symmetries present in the code generation and accusation algorithm, it is usually assumed that the attackers apply a strategy that acts at every location independently. Furthermore, we assume that the colluders take equal risks. In such an attack model, the colluders’ decision in location i depends only on the tallies $M_{i,\alpha} = |\{j \in \mathcal{C} | X_{i,j} = \alpha\}|$ (with $\alpha \in \mathcal{A}$). The tallies satisfy $|\mathbf{M}_i| = c$, and they are multinomial-distributed, $f_{\mathbf{m}|\mathbf{p}} = \binom{c}{\mathbf{m}} \mathbf{p}^m$. The attack strategy may be probabilistic.

C. Combined-Digit Model (CDM)

The CDM [16] allows colluders to mix symbols and to introduce noise (see Figure 1). In each location, the symbols that are mixed are assumed to have equal power. The set of symbols that the colluders choose to mix is denoted as $\Psi \subseteq \mathcal{A}$ with $m_\alpha > 0$ for each $\alpha \in \Psi$. The attack strategy is parametrized by a set of probabilities $f_{\psi|\mathbf{m}}$. The tracer has a detector that outputs a set $\Phi \subseteq \mathcal{A}$ of observed symbols. The joint effects of the noise and the mixing lead to probability distributions $f_{\Phi|\Psi}$, where it is possible that the noise introduces symbols in Φ that are absent in Ψ . Simple-decoder score systems were introduced in [16], [17].

The CDM reduces to the RDM when the noise strength is sent to zero and the detector unerringly observes $\Phi = \Psi$,

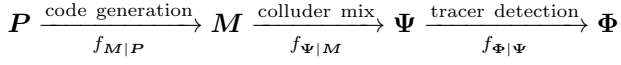


Fig. 1. A schematic depiction of the CDM.

forcing the colluders to output a single symbol, $\Psi = \{Y\}$. For the RDM, a strategy is parametrized by a set of probabilities $f_{y|m}$.

D. Performance; moments of the scores

The performance of bias-based tracing schemes can for a large part be characterized by looking merely at the first and second moment of the innocent and guilty scores. (This holds especially at large c , where the large code length induces an almost-Gaussian shape of the score probability distributions.)

For an innocent user j , we define the mean and variance as

$$\tilde{\mu}_{\text{inn}} := \mathbb{E}[S_{i,j}] \quad (4)$$

$$\tilde{\sigma}_{\text{inn}}^2 := \text{Var}[S_{i,j}] = \mathbb{E}[(S_{i,j} - \tilde{\mu}_{\text{inn}})^2] = \mathbb{E}[S_{i,j}^2] - \tilde{\mu}_{\text{inn}}^2, \quad (5)$$

where the index $i \in [\ell]$ is arbitrary. The expectation \mathbb{E} is taken over the random variables P_i , $X_{i,j}$, and Y_i (in the CDM Ψ_i and Φ_i instead of Y_i). We call a suspicion function centered if it yields $\tilde{\mu}_{\text{inn}} = 0$ and normalized if $\tilde{\sigma}_{\text{inn}}^2 = 1$. For the coalition we define $S_{i,c} := \sum_{j \in \mathcal{C}} S_{i,j}$. The moments are

$$\tilde{\mu}_c := \mathbb{E}[S_{i,c}] \quad (6)$$

$$\tilde{\sigma}_c^2 := \text{Var}[S_{i,c}] = \mathbb{E}[(S_{i,c} - \tilde{\mu}_c)^2] = \mathbb{E}[S_{i,c}^2] - \tilde{\mu}_c^2 \quad (7)$$

again with arbitrary index i . If the Gaussian approximation holds, then the sufficient code length is proportional to $(\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}})^{-2}c^2$ [7]. We will use the fraction $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}}$ as a performance indicator.

III. OPTIMAL SUSPICION FUNCTIONS

We consider suspicion functions h other than the function g given in (3). We derive suspicion functions that maximize the performance indicator $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}}$, in the CDM as well as the RDM. Without loss of generality, we will consider only suspicion functions that are centered ($\tilde{\mu}_{\text{inn}} = 0$) and normalized ($\tilde{\sigma}_{\text{inn}} = 1$). We use the standard approach of Lagrange functionals; we use constraint multipliers $\lambda_1, \lambda_2 \in \mathbb{R}$ to enforce the constraints $\tilde{\mu}_{\text{inn}} = 0$ and $\tilde{\sigma}_{\text{inn}} = 1$. We define the functional

$$L(h, \lambda_1, \lambda_2) = \tilde{\mu}_c - \lambda_1 \tilde{\mu}_{\text{inn}} - \frac{1}{2} \lambda_2 (\tilde{\sigma}_{\text{inn}}^2 - 1), \quad (8)$$

where $\tilde{\mu}_{\text{inn}}$, $\tilde{\sigma}_{\text{inn}}$ and $\tilde{\mu}_c$ depend on the function h as specified in (4-6). The optimal h is found by solving the set of equations $\delta L/\delta h = 0$, $\partial L/\partial \lambda_1 = 0$ and $\partial L/\partial \lambda_2 = 0$. The solution depends on the arguments of h : in the CDM the sub-score of user j in location i is typically a function of $X_{i,j}$, Φ_i and P_i ; in the RDM a function of $X_{i,j}$, Y_i and P_i .

A. Optimal Suspicion Functions in the Combined-Digit Model

We present a number of lemmas leading up to the main theorem of this section, which shows the solution obtained by the Lagrangian approach. The conditional probabilities that appear in the lemmas are related as follows:

$f_{\psi|p} = \sum_{\mathbf{m}} f_{\psi|m} f_{\mathbf{m}|p}$ and $f_{\phi|p} = \sum_{\psi} f_{\phi|\psi} f_{\psi|p}$. The numbers $f_{\phi|\psi}$ are fixed parameters of the CDM independent of the strategy.

Lemma 1. *An optimal suspicion function of the form $h(x, \phi, \psi, p)$ does not depend on ϕ . An optimal suspicion function of the form $h(x, \phi, \psi, \mathbf{m}, p)$ depends neither on ϕ nor ψ .*

Proof sketch: The set ψ contains more information about the attackers than the set ϕ . Likewise, the tallies \mathbf{m} contain more information than ψ . ■

We will give the full proof after Theorem 1.

To determine the optimal suspicion functions of the increasingly general forms $h(x, \phi, p)$, $h(x, \phi, \psi, p)$, and $h(x, \phi, \psi, \mathbf{m}, p)$, it suffices to study the forms $h_{\Phi}(x, \phi, p)$, $h_{\Psi}(x, \psi, p)$, and $h_M(x, \mathbf{m}, p)$, respectively.

Lemma 2. *Let h be of the form $h_{\Phi}(x, \phi, p)$ and define*

$$T_{\Phi}(x, \phi, p) := \frac{\mathbb{E}_{M|P}[M_X f_{\phi|M}]}{c p_X f_{\phi|P}} = \frac{1}{c} \frac{\partial \ln f_{\phi|P}}{\partial p_X} \Big|_{|P|=1} + 1. \quad (9)$$

Then $\tilde{\mu}_c = c \cdot \mathbb{E}[T_{\Phi} h]$ and $\mathbb{E}[T_{\Phi}] = 1$.

Proof: We write (6) as

$$\tilde{\mu}_c = \mathbb{E}_P \mathbb{E}_{M|P} \mathbb{E}_{\Phi|M} \sum_{x \in \mathcal{A}} M_x h(x, \Phi, P) \quad (10)$$

$$= \mathbb{E}_P \mathbb{E}_{M|P} \mathbb{E}_{\Phi|P} \frac{f_{\Phi|M}}{f_{\Phi|P}} \mathbb{E}_{X|P} \frac{M_X}{P_X} h(X, \Phi, P) \quad (11)$$

$$= \mathbb{E}_P \mathbb{E}_{\Phi|P} \mathbb{E}_{X|P} \left[\frac{\mathbb{E}_{M|P}[M_X f_{\Phi|M}]}{P_X f_{\Phi|P}} h(X, \Phi, P) \right] \quad (12)$$

$$= c \mathbb{E}[T \cdot h]. \quad (13)$$

Furthermore, $\mathbb{E}_{X|P}[m_X] = c p_X$ and $f_{\phi|p} = \mathbb{E}_{M|P}[f_{\phi|M}]$, so

$$\mathbb{E}_{X|P}[T] = \mathbb{E}_{X|P} \left[\frac{\mathbb{E}_{M|P}[M_X f_{\phi|M}]}{c p_X f_{\phi|P}} \right] = 1. \quad (14)$$

To be able to take the partial derivative $\frac{\partial \ln f_{\phi|P}}{\partial p_X}$, the components p_0, \dots, p_{q-1} are assumed to be functionally independent. In particular, we do not assume $|p| = 1$ during differentiation. Since $f_{\mathbf{m}|p} = \frac{1}{|p|^c} \binom{c}{\mathbf{m}} p^{\mathbf{m}}$, we find

$$f_{\phi|p} = \mathbb{E}_{M|P}[f_{\phi|M}] = \frac{1}{|p|^c} \sum_{\mathbf{m}} \binom{c}{\mathbf{m}} p^{\mathbf{m}} f_{\phi|\mathbf{m}}. \quad (15)$$

$$\frac{\partial \ln f_{\phi|P}}{\partial p_X} = \frac{\frac{1}{p_X} \sum_{\mathbf{m}} \binom{c}{\mathbf{m}} p^{\mathbf{m}} m_X f_{\phi|\mathbf{m}}}{\sum_{\mathbf{m}} \binom{c}{\mathbf{m}} p^{\mathbf{m}} f_{\phi|\mathbf{m}}} - \frac{c}{|p|} \quad (16)$$

$$= \frac{\mathbb{E}_{M|P}[M_X f_{\phi|M}]}{p_X f_{\phi|P}} - \frac{c}{|p|}. \quad (17)$$

So $\frac{1}{c} \frac{\partial \ln f_{\phi|P}}{\partial p_X} \Big|_{|p|=1} + 1 = T_{\Phi}(x, \phi, p)$. ■

When the colluders output is known the optimal suspicion function is derived as follows:

Lemma 3. Let h be of the form $h_{\Psi}(x, \psi, \mathbf{p})$ and define

$$T_{\Psi}(x, \psi, \mathbf{p}) := \frac{\mathbb{E}_{\mathbf{M}|\mathbf{P}}[M_x f_{\psi|\mathbf{M}}]}{c p_x f_{\psi|\mathbf{P}}} = \frac{1}{c} \frac{\partial \ln f_{\psi|\mathbf{P}}}{\partial p_x} \Big|_{|\mathbf{p}|=1} + 1. \quad (18)$$

Then $\tilde{\mu}_c = c \cdot \mathbb{E}[T_{\Psi}h]$ and $\mathbb{E}[T_{\Psi}] = 1$.

Proof: We write (6) as

$$\tilde{\mu}_c = \mathbb{E}_{\mathbf{P}} \mathbb{E}_{\mathbf{M}|\mathbf{P}} \mathbb{E}_{\Psi|\mathbf{M}} \sum_{x \in \mathcal{A}} M_x h(x, \Psi, \mathbf{P}). \quad (19)$$

Note the similarity between (19) and (10). The proof proceeds analogously with Ψ instead of Φ . ■

When even the tallies of the coalition symbols are known the optimal suspicion function is derived as follows:

Lemma 4. Let h be of the form $h_{\mathbf{M}}(x, \mathbf{m}, \mathbf{p})$ and define

$$T_{\mathbf{M}}(x, \mathbf{m}, \mathbf{p}) := \frac{m_x}{c p_x} = \frac{1}{c} \frac{\partial \ln f_{\mathbf{m}|\mathbf{P}}}{\partial p_x} \Big|_{|\mathbf{p}|=1} + 1. \quad (20)$$

Then $\tilde{\mu}_c = c \cdot \mathbb{E}[T_{\mathbf{M}}h]$, $\mathbb{E}[T_{\mathbf{M}}] = 1$, and $\text{Var}[T_{\mathbf{M}}] = \frac{q-1}{c}$.

Proof: We write (6) as

$$\tilde{\mu}_c = \mathbb{E}_{\mathbf{P}} \mathbb{E}_{\mathbf{M}|\mathbf{P}} \sum_{x \in \mathcal{A}} M_x h(x, \mathbf{M}, \mathbf{P}) \quad (21)$$

$$= \mathbb{E}_{\mathbf{P}} \mathbb{E}_{\mathbf{M}|\mathbf{P}} \mathbb{E}_{X|\mathbf{P}} \left[\frac{M_X}{P_X} h(X, \mathbf{M}, \mathbf{P}) \right] = c \mathbb{E}[T \cdot h]. \quad (22)$$

Furthermore, $\mathbb{E}_{\mathbf{M}|\mathbf{P}}[M_x] = c p_x$, so

$$\mathbb{E}[T] = \mathbb{E}_{\mathbf{P}} \mathbb{E}_{\mathbf{M}|\mathbf{P}} \mathbb{E}_{X|\mathbf{P}} \left[\frac{M_X}{c P_X} \right] = \mathbb{E}_{\mathbf{P}} \mathbb{E}_{X|\mathbf{P}}[1] = 1. \quad (23)$$

Also

$$\text{Var}[T] = \mathbb{E}_{\mathbf{P}} \mathbb{E}_{X|\mathbf{P}} \mathbb{E}_{\mathbf{M}|\mathbf{P}} \left(\frac{M_X}{c P_X} - 1 \right)^2 \quad (24)$$

$$= \mathbb{E}_{\mathbf{P}} \mathbb{E}_{X|\mathbf{P}} \text{Var}_{\mathbf{M}|\mathbf{P}} \left[\frac{M_X}{c P_X} \right] = \mathbb{E}_{\mathbf{P}} \mathbb{E}_{X|\mathbf{P}} \left[\frac{c p_X (1 - p_X)}{c^2 p_X^2} \right] \quad (25)$$

$$= \frac{1}{c} \mathbb{E}_{\mathbf{P}} \sum_{x \in \mathcal{A}} (1 - p_x) = (q - 1)/c. \quad (26)$$

Also, $\frac{\partial \ln f_{\mathbf{m}|\mathbf{P}}}{\partial p_x} = \frac{m_x}{p_x} - \frac{c}{|\mathbf{p}|}$ and thus $\frac{1}{c} \frac{\partial \ln f_{\mathbf{m}|\mathbf{P}}}{\partial p_x} \Big|_{|\mathbf{p}|=1} + 1 = T_{\mathbf{M}}(x, \mathbf{m}, \mathbf{p})$. ■

Theorem 1. In each of the cases above, the centered and normalized suspicion function that maximizes $\tilde{\mu}_c$ is

$$h = (T - \mathbb{E}[T]) / \sqrt{\text{Var}[T]} \quad (27)$$

and the expected coalition score is $\tilde{\mu}_c = c \cdot \sqrt{\text{Var}[T]}$.

Proof: Define the Lagrangian

$$L(h, \lambda_1, \lambda_2) := c \mathbb{E}[Th] - \lambda_1 \mathbb{E}[h] - \frac{1}{2} \lambda_2 (\mathbb{E}[h^2] - 1) \quad (28)$$

with the two Lagrange multipliers λ_1 and λ_2 enforcing that the function is centered and normalized respectively. Let h be such that $\frac{\delta L}{\delta h} = 0$. Then $D(cT - \lambda_1 - \lambda_2 h) = 0$ (where D is the

product of the probability densities of the random variables), i.e. $h = \frac{cT - \lambda_1}{\lambda_2}$. The first constraint, $\tilde{\mu}_{\text{inn}} = 0$, implies that $\lambda_1 = c \mathbb{E}[T]$ and the second constraint, $\tilde{\sigma}_{\text{inn}}^2 = 1$, implies that $\lambda_2^2 = \mathbb{E}(cT - \lambda_1)^2 = c^2 \text{Var}[T]$. ■

Now that we have seen the proof technique for Theorem 1, we can state the full proof of Lemma 1:

Proof of Lemma 1: To determine the optimal suspicion function of the form $h(x, \psi, \mathbf{p})$ in the proof of Theorem 1 we defined the Lagrangian

$$L(h, \lambda_1, \lambda_2) := c \mathbb{E}[Th] - \lambda_1 \mathbb{E}[h] - \frac{1}{2} \lambda_2 (\mathbb{E}[h^2] - 1). \quad (29)$$

where $\mathbb{E}[\dots] = \mathbb{E}_{\mathbf{P}} \mathbb{E}_{\Psi|\mathbf{P}} \mathbb{E}_{X|\mathbf{P}}[\dots]$. The Euler-Lagrange equation was $D(cT - \lambda_1 - \lambda_2 h) = 0$ with $D = f_{\mathbf{P}} f_{\psi|\mathbf{P}} f_{X|\mathbf{P}}$.

Instead, to determine the optimal suspicion function of the form $h(x, \phi, \psi, \mathbf{p})$, we would define the same Lagrangian, but now with $\mathbb{E}[\dots] = \mathbb{E}_{\mathbf{P}} \mathbb{E}_{\Psi|\mathbf{P}} \mathbb{E}_{\Phi|\Psi} \mathbb{E}_{X|\mathbf{P}}[\dots]$. We obtain the same Euler-Lagrange equation but now with $D = f_{\mathbf{P}} f_{\psi|\mathbf{P}} f_{\phi|\psi} f_{X|\mathbf{P}}$.

In both cases, we draw the same conclusion: that $cT - \lambda_1 - \lambda_2 h = 0$. We therefore find that the optimal suspicion function of the form $h(x, \phi, \psi, \mathbf{p})$ is the one we found in Lemma 3 of the form $h(x, \psi, \mathbf{p})$.

Likewise, the optimal suspicion function of the form $h(x, \phi, \psi, \mathbf{m}, \mathbf{p})$ is the one we found in Lemma 4 of the form $h(x, \mathbf{m}, \mathbf{p})$. ■

Our suspicion functions have a close relation with Neyman-Pearson scores, as shown in the following proposition.

Proposition 5. In all three cases $T_{\Phi}(x, \phi, \mathbf{p})$, $T_{\Psi}(x, \psi, \mathbf{p})$ and $T_{\mathbf{M}}(x, \mathbf{m}, \mathbf{p})$ for the function T it holds that

$$T(x, \square, \mathbf{p}) \propto \frac{\mathbb{P}[j \in \mathcal{C} | x, \square, \mathbf{p}]}{\mathbb{P}[j \notin \mathcal{C} | x, \square, \mathbf{p}]}, \quad (30)$$

and thus T is a Neyman-Pearson score.

Proof: The Neyman-Pearson score for testing a hypothesis H given evidence e is given by the likelihood ratio $\mathbb{P}[H = \text{True}|e] / \mathbb{P}[H = \text{False}|e]$. Our hypothesis is $H = (j \in \mathcal{C})$ for a user $j \in [n]$, and we consider the evidence $e = (x, \phi, \mathbf{p})$ available in one location. (The proof for all the other cases is analogous.) Then the Neyman-Pearson score is

$$\frac{\mathbb{P}[j \in \mathcal{C} | x, \phi, \mathbf{p}]}{\mathbb{P}[j \notin \mathcal{C} | x, \phi, \mathbf{p}]} = \frac{\mathbb{P}[j \in \mathcal{C}, x, \phi, \mathbf{p}]}{\mathbb{P}[j \notin \mathcal{C}, x, \phi, \mathbf{p}]} \quad (31)$$

$$= \frac{\mathbb{P}[j \in \mathcal{C}] f_{\mathbf{P}} f_{x|\mathbf{P}} f_{\phi|x, \mathbf{p}, j \in \mathcal{C}}}{\mathbb{P}[j \notin \mathcal{C}] f_{\mathbf{P}} f_{x|\mathbf{P}} f_{\phi|\mathbf{P}}} \quad (32)$$

$$\propto \frac{f_{\phi|x, \mathbf{p}, j \in \mathcal{C}}}{f_{\phi|\mathbf{P}}} \quad (33)$$

$$= \frac{1}{f_{\phi|\mathbf{P}}} \sum_{\mathbf{m}: m_x \geq 1} \binom{c-1}{\mathbf{m} - \mathbf{e}_x} p^{m - \mathbf{e}_x} f_{\phi|\mathbf{m}} \quad (34)$$

$$= \frac{1}{f_{\phi|\mathbf{P}}} \sum_{\mathbf{m}} \frac{m_x}{c p_x} \binom{c}{\mathbf{m}} p^{\mathbf{m}} f_{\phi|\mathbf{m}} \quad (35)$$

$$= \frac{1}{f_{\phi|\mathbf{p}}} \mathbb{E}_{M|\mathbf{p}} \left[\frac{M_x}{cp_x} f_{\phi|M} \right]. \quad (36)$$

Here e_x is a length q vector containing a 1 in position x and zero elsewhere. The a priori probability $\mathbb{P}[j \in \mathcal{C}]$ is a constant. It is equal for all users if the tracer has no prior knowledge about the coalition. ■

Several things are worth noting about these results.

- (i) In the proof of Theorem 1 it is not necessary to specify the bias distribution. Though $\tilde{\mu}_{\mathcal{C}}$ is a functional of both h and $f_{\mathcal{P}}$, the optimization of h does not depend on $f_{\mathcal{P}}$.
- (ii) In all three cases the result for h depends on information that the tracer usually does not have. (The strategy $f_{\psi|m}$ in Lemmas 2 and 3; the tallies \mathbf{m} in Lemma 4.) When a function h_{Φ} , for some guessed strategy, is used to compute scores, there is no guarantee that the attackers are actually adhering to that guessed strategy. Such ‘mismatched’ situations will be discussed (for the RDM) in the remainder of this paper.
- (iii) We can think of two ways in which the \mathbf{m} -dependent result of Lemma 4, $h(x, y, \mathbf{p}) = \left(\frac{m_x}{cp_x} - 1\right) \sqrt{\frac{c}{q-1}}$, can be used in practice. First, it could be employed in the EM algorithm [12]. The EM procedure estimates a strategy based on the symbols received by the candidate coalition, and then uses this estimate to adapt the suspicion function. Our h function could be used to directly assign scores to all users, *skipping the strategy estimation step*. This would speed up each iteration of the EM algorithm and avoid the statistical inaccuracies in the estimation. (Of course, inaccuracies due to a wrongly guessed coalition remain, and may even increase.) Secondly, this h function can be used as a consistency check in the following way. Suppose that, by some means, a candidate coalition $\hat{\mathcal{C}}$ has been tentatively identified. Then one computes a score $\left(\frac{m_x}{cp_x} - 1\right) \sqrt{\frac{c}{q-1}}$ for all users, where the tally m_x is based on $\hat{\mathcal{C}}$ and the user’s symbol x . If $\hat{\mathcal{C}}$ equals the actual coalition, one should see a huge score difference between innocent users and the colluders. Exploration of these ideas is left for future work.
- (iv) The expression $\partial \ln f / \partial p_x$ in all three cases has the form of a Fisher score, being the derivative of the logarithm of a conditional probability with respect to the conditioning variable. We suspect that this form is no coincidence. However, the intuitive meaning of the associated ‘game’ (guessing \mathbf{p} from y) is not immediately obvious. Asymptotically \mathbf{m} tends to $c\mathbf{p}$. We hypothesize that the game ‘guess \mathbf{p} from y ’ is asymptotically equivalent to ‘guess \mathbf{m} from y ’. The latter is a known formulation of the tracing problem.
- (v) Our result in Proposition 5 is different from the Neyman-Pearson score in [13], where the whole sequence $(Y_i)_{i \in [\ell]}$ was considered.

B. Optimal Suspicion Functions in the Restricted-Digit Model

The Restricted-Digit Model is a special case of the Combined-Digit Model.

Corollary 6. *Let h be of the form $h_Y(x, y, \mathbf{p})$ and define*

$$T_Y(x, y, \mathbf{p}) := \frac{\mathbb{E}_{M|\mathbf{p}}[M_x f_{y|M}]}{cp_x f_{y|\mathbf{p}}} = \frac{1}{c} \frac{\partial \ln f_{y|\mathbf{p}}}{\partial p_x} \Big|_{|\mathbf{p}|=1} + 1. \quad (37)$$

Then $\tilde{\mu}_{\mathcal{C}} = c \cdot \mathbb{E}[T_Y h]$ and $\mathbb{E}[T_Y] = 1$.

Proof: The optimal h function in the RDM case follows straightforwardly from Lemma 2 and Theorem 1 by taking the limit of zero noise and perfect detection of all mixed symbols, leading to $\Phi = \Psi = \{Y\}$, with $Y \in \mathcal{A}$. ■

In the RDM, Lemma 4 and Theorem 1 hold without change. Note that the Marking Assumption is not invoked to obtain Corollary 6. Hence Corollary 6 is valid in a more general setting, as long as the colluders produce a single symbol which is unerringly detected by the tracer.

Note also that (37) with $q = 2$ matches the expression given by Charpentier et al. [12] (which only considered the binary case).

C. Strongly Centered and Normalized Suspicion Functions

In (8) we required our optimal score functions to be centered and normalized. The normalization was done without loss of generality, since scores can be rescaled arbitrarily. The symmetric Tardos suspicion function was chosen to satisfy stronger properties: it is both centered and normalized, no matter what the pirate symbol y or the bias vector \mathbf{p} are (and no matter what the attack strategy or the bias distribution is for that matter). These properties are captured in the following definition.

We call a suspicion function $h(x, y, \mathbf{p})$ *strongly centered* if $\mathbb{E}_{X|\mathbf{p}}[h(X, y, \mathbf{p})] = 0$ and *strongly normalized* if $\mathbb{E}_{X|\mathbf{p}}[h^2(X, y, \mathbf{p})] = 1$.

We show that even when the score function does not match the pirate strategy, the optimal score functions derived in the previous section remain centered but not necessarily normalized.

Lemma 7. *Each optimal suspicion function (see Theorem 1) is strongly centered. So is the symmetric Tardos function.*

Proof: This follows directly from (14). ■

If we wanted to find optimal suspicion functions that are both *strongly centered* and *strongly normalized*, like the symmetric Tardos suspicion function, in (8) we should require $\mathbb{E}_{X|\mathbf{p}}[h(X, y, \mathbf{p})] = 0$ and $\text{Var}_{X|\mathbf{p}}[h(X, y, \mathbf{p})] = 1$. Since our optimal suspicion functions already turned out to be strongly centered, in Theorem 1 only the normalizing constant changes:

Corollary 8. *The strongly centered and strongly normalized suspicion function that maximizes $\tilde{\mu}_{\mathcal{C}}$ is*

$$h = (T - \mathbb{E}_{X|\mathbf{p}}[T]) / \sqrt{\text{Var}_{X|\mathbf{p}}[T]} \quad (38)$$

and the expected coalition score is $\tilde{\mu}_c = c \cdot \sqrt{\text{Var}_{X|\mathbf{p}}[T]}$.

Proof: Define the Lagrangian

$$\begin{aligned} L(h, \lambda_1, \lambda_2) := & c \mathbb{E}_{X|\mathbf{p}}[T(X, y, \mathbf{p})h(X, y, \mathbf{p})] + \\ & - \lambda_1 \mathbb{E}_{X|\mathbf{p}}[h(X, y, \mathbf{p})] + \\ & - \frac{1}{2} \lambda_2 (\mathbb{E}_{X|\mathbf{p}}[h^2(X, y, \mathbf{p})] - 1) \end{aligned} \quad (39)$$

with the two Lagrange multipliers λ_1 and λ_2 enforcing that the function is centered and normalized respectively. Let h be such that $\frac{\delta L(h, \lambda_1, \lambda_2)}{\delta h(x, y, \mathbf{p})} = 0$. Then

$$f_{x|\mathbf{p}}(cT(x, y, \mathbf{p}) - \lambda_1 - \lambda_2 h(x, y, \mathbf{p})) = 0, \quad (40)$$

i.e. $h(x, y, \mathbf{p}) = \frac{cT(x, y, \mathbf{p}) - \lambda_1}{\lambda_2}$. The first constraint that $h(x, y, \mathbf{p})$ is strongly centered implies that $\lambda_1 = c \mathbb{E}_{X|\mathbf{p}}[T(x, y, \mathbf{p})]$ and the second constraint that $h(x, y, \mathbf{p})$ is strongly normalized implies that

$$\lambda_2^2 = \mathbb{E}_{X|\mathbf{p}}(cT(x, y, \mathbf{p}) - \lambda_1)^2 = c^2 \text{Var}_{X|\mathbf{p}}[T(x, y, \mathbf{p})]. \quad (41)$$

■

D. Building a Traitor Tracing Scheme

Now that we have described our new optimal suspicion function, there is one caveat left to address. As noted before, when a suspicion function h , for some guessed strategy, is used to compute scores, there is no guarantee that the attackers are actually adhering to that guessed strategy. In particular, this means that we no longer have the property $\tilde{\sigma}_{\text{inn}}^2 = 1$ which the Tardos suspicion function enjoys.

As a result, we can not simply plug our suspicion function into the Tardos traitor tracing scheme, since such a scheme typically accuses a user when he exceeds a fixed threshold. Traditionally, the use of a fixed threshold is possible since the scaling of the scores is taken care of by the property $\tilde{\sigma}_{\text{inn}}^2 = 1$. Thus, ideally, we would like to have a normalized suspicion function. This can be achieved by scaling all scores (i.e. scaling the function h) by a factor $\tilde{\sigma}_{\text{inn}}$. Unfortunately, we have no way of knowing $\tilde{\sigma}_{\text{inn}}$ when the exact collusion strategy is unknown. However, in practice, when the number of colluders is not too large compared to the number of users, we can estimate $\tilde{\sigma}_{\text{inn}}$ by the *observed* $\tilde{\sigma}$ of the complete population. If we scale all scores by a factor $\tilde{\sigma}$ by replacing them with $S_j/\tilde{\sigma}$ for every user j , the scheme will perform well against any collusion strategy.

IV. DEFENDING AGAINST COMMON COLLUSION STRATEGIES

From this point onward, we consider only the RDM. For a number of often-studied strategies we compute the optimal suspicion function. We investigate the situation where the actual attack is indeed the one for which the h -function was designed (a ‘‘match’’), as well as mismatches. We will call the ‘‘optimal suspicion function against strategy A’’ the A-defense. The following sections will focus on defenses against five often-considered strategies. In short, these strategies can be described as follows:

- 1) *Interleaving attack (Section V)*: The interleaving attack randomly selects an attacker and outputs his symbol.
- 2) *All-high (all-1) attack (Section VI)*: The all-high attack is special as it breaks the symbol symmetry. It assumes that the alphabet can be ordered in some meaningful way, and outputs the largest received symbol. In the binary case $q = 2$ this attack is known as the all-1 attack, as it will output a 1 if the coalition has received one.
- 3) *Random-symbol (coin-flip) attack (Section VII)*: The random-symbol attack randomly selects a received symbol, irrespective of the tally vector \mathbf{m} , and outputs it. In the binary case $q = 2$ this attack is known as the coin-flip attack.
- 4) *Majority voting attack (Section VIII)*: The majority voting attack outputs the symbol that was received most often by the coalition. In case multiple symbols are received equally often, a random symbol is chosen among them.
- 5) *Minority voting attack (Section IX)*: The minority voting attack outputs the symbol that was received least often (but at least once) by the coalition. When multiple symbols are received equally often, a random symbol is chosen among them.

A detailed description of each attack will be given at the start of its section. We also dedicate a section (Section X) to analyzing the performance of the traditional symmetrized Tardos suspicion function against these attacks.

V. INTERLEAVING DEFENSE

A. Optimal defense

The interleaving attack $f_{y|\mathbf{m}} = m_y/c$ randomly selects an attacker and outputs his symbol.

Proposition 9. *Against the interleaving attack, the quantity T is given by $T(x, y, \mathbf{p}) = 1 + \frac{1}{c}(\delta_{x,y}/p_y - 1)$, and the optimal suspicion function is*

$$h(x, y, \mathbf{p}) = \frac{1}{\sqrt{q-1}} \left(\frac{\delta_{x,y}}{p_y} - 1 \right). \quad (42)$$

Proof: We find

$$f_{y|\mathbf{p}} = \frac{1}{c|\mathbf{p}|^c} \sum_{\mathbf{m}} \binom{c}{\mathbf{m}} \mathbf{p}^{\mathbf{m}} m_y = \frac{p_y}{c|\mathbf{p}|^c} \frac{\partial |\mathbf{p}|^c}{\partial p_y} = \frac{p_y}{|\mathbf{p}|}. \quad (43)$$

Thus

$$\left. \frac{\partial \ln f_{y|\mathbf{p}}}{\partial p_x} \right|_{|\mathbf{p}|=1} = \frac{\delta_{x,y}}{p_y} - 1, \quad (44)$$

so $T(x, y, \mathbf{p}) = 1 + \frac{1}{c}(\delta_{x,y}/p_y - 1)$. Also,

$$\text{Var}[T] = \mathbb{E}(T - 1)^2 \quad (45)$$

$$= \frac{1}{c^2} \mathbb{E}_{\mathbf{P}} \mathbb{E}_{Y|\mathbf{P}} \mathbb{E}_{X|\mathbf{P}} \left[(\delta_{x,y}/p_y - 1)^2 \right] \quad (46)$$

$$= \frac{1}{c^2} \mathbb{E}_{\mathbf{P}} \mathbb{E}_{Y|\mathbf{P}} \left[P_Y \left(\frac{1 - P_Y}{P_Y} \right)^2 + \sum_{x \neq y} P_x \right] \quad (47)$$

$$= \frac{1}{c^2} \mathbb{E}_{\mathbf{P}} \mathbb{E}_{Y|\mathbf{P}} \left[\frac{1 - P_Y}{P_Y} \right] \quad (48)$$

$$= \frac{1}{c^2} \mathbb{E}_{\mathbf{P}} [q - 1] = \frac{q - 1}{c^2}. \quad (49)$$

■ and, with (118),

$$\mathbb{E}_{\mathcal{P}}\mathbb{E}_{Y|\mathcal{P}}\left[\frac{1}{P_Y}\right] = \sum_{y=0}^{q-1}\mathbb{E}_{\mathcal{P}}\left[\frac{A_{y+1}^c - A_y^c}{P_y}\right]. \quad (59)$$

The performance of the interleaving defense against the interleaving attack is given in the following lemma.

Proposition 10. *When the interleaving attack is used against the interleaving defense, then $\tilde{\mu}_c = \sqrt{q-1}$, achieving capacity for any $f_{\mathcal{P}}$.*

Proof: From Theorem 1 we know that $\tilde{\mu}_c = c \cdot \sqrt{\text{Var}[T]}$. Combining this with (49) yields $\tilde{\mu}_c = \sqrt{q-1}$, which corresponds to the capacity [19]. ■

When $x = y$, the h is positive and increasing in p_y (rare events raise more suspicion). When $x \neq y$, it is negative and constant, in contrast to (3). The h is independent of c .

Lemma 11. *If the tracer uses the interleaving defense, then, no matter what attack is used,*

$$\tilde{\mu}_c = \frac{c}{\sqrt{q-1}} \left(-1 + \mathbb{E}_{\mathcal{P}}\mathbb{E}_{Y|\mathcal{P}}[T(Y, Y, \mathcal{P})]\right) \quad (50)$$

and

$$\tilde{\sigma}_{\text{inn}}^2 = \frac{1}{q-1} \left(-1 + \mathbb{E}_{\mathcal{P}}\mathbb{E}_{Y|\mathcal{P}}\left[\frac{1}{P_Y}\right]\right). \quad (51)$$

where T belongs to the attack.

Proof: Using the interleaving defense from (42), we find

$$\tilde{\mu}_c = \mathbb{E}[T \cdot h] = \frac{c}{\sqrt{q-1}} \left(\mathbb{E}\left[T(X, Y, \mathcal{P}) \frac{\delta_{X,Y}}{P_Y}\right] - 1\right). \quad (52)$$

Also

$$h^2(x, y, \mathcal{P}) = \frac{1}{q-1} \left(\frac{\delta_{x,y}}{p_y} \left(\frac{1}{p_y} - 2\right) + 1\right), \quad (53)$$

$$\text{so } \tilde{\sigma}_{\text{inn}}^2 = \mathbb{E}[h^2] \quad (54)$$

$$= \frac{1}{q-1} \left(1 + \mathbb{E}_{\mathcal{P}}\mathbb{E}_{Y|\mathcal{P}}\left[\frac{1}{P_Y} - 2\right]\right). \quad (55)$$

We can explicitly calculate the performance against the all-high attack (which is formalized in Proposition 19):

Proposition 12. *If the tracer uses the interleaving defense, but the coalition uses the all-high attack, then*

$$\tilde{\mu}_c = \frac{c}{\sqrt{q-1}} \sum_{y=0}^{q-2} \mathbb{E}_{\mathcal{P}}[A_{y+1}^{c-1}], \text{ and} \quad (56)$$

$$\tilde{\sigma}_{\text{inn}}^2 = \frac{1}{q-1} \left(-1 + \sum_{y=0}^{q-1} \mathbb{E}_{\mathcal{P}}\left[\frac{A_{y+1}^c - A_y^c}{P_y}\right]\right). \quad (57)$$

Proof: Using Lemma 11 with (113), we find

$$\mathbb{E}_{\mathcal{P}}\mathbb{E}_{Y|\mathcal{P}}[T(Y, Y, \mathcal{P})] = \sum_{y=0}^{q-1} \mathbb{E}_{\mathcal{P}}[A_{y+1}^{c-1}] = \sum_{y=0}^{q-2} \mathbb{E}_{\mathcal{P}}[A_{y+1}^c] + 1. \quad (58)$$

If the Dirichlet distribution is used $\tilde{\mu}_c$ will scale as $c^{1-\kappa}$ for large coalitions:

Proposition 13. *Let $f_{\mathcal{P}}$ be the symmetric Dirichlet distribution with cutoff $\delta = 0$. If the tracer uses the interleaving defense, but the colluders use the all-high attack, then*

$$\tilde{\mu}_c = \frac{\Gamma(q\kappa)}{\Gamma((q-1)\kappa)} \frac{c^{1-\kappa}}{\sqrt{q-1}} [1 + \mathcal{O}(1/c)]. \quad (60)$$

Proof: Lemma 11 gives

$\tilde{\mu}_c = \frac{c}{\sqrt{q-1}} (\mathbb{E}_{\mathcal{P}}\mathbb{E}_{Y|\mathcal{P}}[T(Y, Y, \mathcal{P})] - 1)$. Next, using Proposition 19 we get $\tilde{\mu}_c = \frac{c}{\sqrt{q-1}} [-1 + \sum_{y=0}^{q-1} \mathbb{E}_{\mathcal{P}} A_y^{c-1}]$ which can be simplified to $\tilde{\mu}_c = \frac{c}{\sqrt{q-1}} \sum_{y=0}^{q-2} \mathbb{E}_{\mathcal{P}} A_y^{c-1}$. The easiest way to evaluate the expectation is by using the marginal distribution of A_y , which is given by $M(a_y) = a_y^{\kappa-1} (1 - a_y)^{[q-y]\kappa-1} / B(y\kappa, [q-y]\kappa)$. (See derivation at the end of this proof.) This yields

$$\tilde{\mu}_c = \frac{c}{\sqrt{q-1}} \sum_{y=0}^{q-2} \frac{B([q-1-y]\kappa, [y+1]\kappa + c - 1)}{B([q-1-y]\kappa, [y+1]\kappa)} \quad (61)$$

$$= \frac{c}{\sqrt{q-1}} \sum_{b=1}^{q-1} \frac{\Gamma(q\kappa)\Gamma(c-1+b\kappa)}{\Gamma(b\kappa)\Gamma(c-1+q\kappa)}. \quad (62)$$

Next we use the property $\Gamma(x+\alpha)/\Gamma(x+\beta) = x^{\alpha-\beta} [1 + \mathcal{O}(1/x)]$ which holds if $x \gg 1$, $a, b \ll x$, and a, b independent of x . (See e.g. Lemma 7 in [6].) This gives

$$\tilde{\mu}_c = \frac{c}{\sqrt{q-1}} \sum_{b=1}^{q-1} \frac{\Gamma(q\kappa)}{\Gamma(b\kappa)} c^{(b-q)\kappa} [1 + \mathcal{O}\left(\frac{1}{c}\right)]. \quad (63)$$

The dominant term is $b = q-1$, yielding (60). The smaller b values in the sum are terms of relative order $1/c$ or smaller.

Finally we derive the marginal distribution $M(a_y)$. We compute $M(a_y) = \mathbb{E}_{\mathcal{P}} \delta(a_y - \sum_{\alpha=0}^{y-1} p_{\alpha})$,

$$M(a_y) = \int_0^1 d^q p \delta(1 - |p|) \frac{p^{\kappa-1}}{B(\kappa \mathbf{1}_q)} \delta(a_y - \sum_{\alpha=0}^{y-1} p_{\alpha}), \quad (64)$$

where $\mathbf{1}_q$ is a vector consisting of q ones and B is the generalized Beta function. We do the following change of integration variables: for $\alpha < y$ we write $p_{\alpha} = a_y t_{\alpha}$ and for $\alpha \geq y$ we write $p_{\alpha} = (1 - a_y) s_{\alpha}$. This gives $\delta(a_y - \sum_{\alpha=0}^{y-1} p_{\alpha}) = a_y^{-1} \delta(1 - |t|)$ and $\delta(1 - |p|) = (1 - a_y)^{-1} \delta(1 - |s|)$. Furthermore, $d^q p p^{\kappa-1} = d^y t d^{q-y} s a_y^{\kappa} (1 - a_y)^{[q-y]\kappa} t^{\kappa-1} s^{\kappa-1}$. Substitution into (64) gives

$$M(a_y) = \frac{a_y^{\kappa-1} (1 - a_y)^{[q-y]\kappa-1}}{B(\kappa \mathbf{1}_q)} \left[\int_0^1 d^y t \delta(1 - |t|) t^{\kappa-1} \right] \cdot \left[\int_0^1 d^{q-y} s \delta(1 - |s|) s^{\kappa-1} \right] \quad (65)$$

$$= \frac{a_y^{\kappa-1} (1 - a_y)^{[q-y]\kappa-1}}{B(\kappa \mathbf{1}_q)} B(\kappa \mathbf{1}_y) B(\kappa \mathbf{1}_{q-y}). \quad (66)$$

Simplification of the Beta functions gives the density $M(a_y)$ as listed earlier in this proof. ■

We now investigate the binary case $q = 2$. We can then rephrase Proposition 12 as

Corollary 14. *Let $q = 2$. If the tracer uses the interleaving defense, but the coalition uses the all-1 attack, then $\tilde{\mu}_C = c \mathbb{E}_{\mathbf{P}} [P_0^{c-1}]$ and $\tilde{\sigma}_{\text{inn}}^2 = -1 + \mathbb{E}_{\mathbf{P}} [P_0^{c-1} + \frac{1-P_0^c}{P_1}]$.*

Proof: $A_1 = P_0$ and $A_2 = P_0 + P_1 = 1$. ■

In the binary case, we obtain explicit results for the coin-flip attack (as formalized in Proposition 29) against the interleaving defense:

Proposition 15. *Let $q = 2$. If the tracer uses the interleaving defense, but the coalition uses the coin-flip attack, then*

$$\tilde{\mu}_C = \frac{1}{2}c \mathbb{E}_{\mathbf{P}} [P_0^{c-1} + P_1^{c-1}] \text{ and} \quad (67)$$

$$\tilde{\sigma}_{\text{inn}}^2 = -1 + \mathbb{E}_{\mathbf{P}} \left[\frac{1 + P_0^c - P_1^c}{2P_0} + \frac{1 + P_1^c - P_0^c}{2P_1} \right]. \quad (68)$$

Proof: Using Lemma 11 with (162), we find

$$\mathbb{E}_{\mathbf{P}} \mathbb{E}_{Y|\mathbf{P}} [T(Y, Y, \mathbf{P})] = \frac{1}{2} \sum_{y \in \mathcal{A}} \mathbb{E}_{\mathbf{P}} [1 + p_y^{c-1}] \quad (69)$$

$$= 1 + \frac{1}{2} \mathbb{E}_{\mathbf{P}} [p_0^{c-1} + p_1^{c-1}]. \quad (70)$$

and, with (166),

$$\mathbb{E}_{\mathbf{P}} \mathbb{E}_{Y|\mathbf{P}} \left[\frac{1}{P_Y} \right] = \frac{1}{2} \sum_{y \in \mathcal{A}} \mathbb{E}_{\mathbf{P}} \frac{1 + p_y^c - p_{1-y}^c}{P_y}. \quad (71)$$

Note the similarity between the coin-flip attack and the all-1 attack. For the Dirichlet distribution, this can be analytically shown:

Proposition 16. *Let $q = 2$ and $f_{\mathbf{P}}$ be the symmetric Dirichlet distribution with parameter $\kappa = \frac{1}{2}$ without cutoff. If the tracer uses the interleaving defense and the coalition uses either the all-1 or the coin-flip attack, then*

$$\tilde{\mu}_C = c \cdot B(\kappa, \kappa + c - 1) / B(\kappa, \kappa) \quad \text{and} \quad (72)$$

$$\tilde{\sigma}_{\text{inn}}^2 = -1 + \frac{c}{1 - \kappa} \frac{\Gamma(2\kappa)}{\Gamma(\kappa)} \frac{\Gamma(c + \kappa - 1)}{\Gamma(c + 2\kappa - 1)} + \frac{1 - 2\kappa}{1 - \kappa}. \quad (73)$$

For large c these behave as $\tilde{\mu}_C \propto c^{1-\kappa}$ and $\tilde{\sigma}_{\text{inn}}^2 \propto c^{1-\kappa}$.

Proof: In the case of the coin-flip attack we have

$$\tilde{\mu}_C = \frac{1}{2}c \mathbb{E}_{\mathbf{P}} [P_0^{c-1} + P_1^{c-1}] = c \mathbb{E}_{\mathbf{P}} [P_0^{c-1}] \quad (74)$$

$$= cB(\kappa, \kappa + c - 1) / B(\kappa, \kappa) \quad (75)$$

since $f_{\mathbf{P}}$ is symbol-symmetric.

Also, $f_{y|\mathbf{P}} = \frac{1}{2} + \frac{1}{2}p_y^c - \frac{1}{2}(1 - p_y)^c$. When the interleaving suspicion function is used, (51) tells us that $\tilde{\sigma}_{\text{inn}}^2 = -1 +$

$\mathbb{E}[1/P_Y]$. We have

$$\mathbb{E} \left[\frac{1}{P_Y} \right] = \sum_{y \in \{0,1\}} \mathbb{E}_{\mathbf{P}} \left[\frac{f_{y|\mathbf{P}}}{P_y} \right] \quad (76)$$

$$= \frac{1}{2} \sum_{y \in \{0,1\}} \mathbb{E}_{\mathbf{P}} \left[\frac{1}{P_y} + P_y^{c-1} - \frac{(1 - P_y)^c}{P_y} \right] \quad (77)$$

$$= \mathbb{E}_{\mathbf{P}} \left[\frac{1}{P_y} + P_y^{c-1} - \frac{(1 - P_y)^c}{P_y} \right] \quad (78)$$

$$= \frac{B(\kappa - 1, \kappa) + B(c + \kappa - 1, \kappa) - B(\kappa - 1, c + \kappa)}{B(\kappa, \kappa)}. \quad (79)$$

In the third line we used the fact that $f_{\mathbf{P}}$ is symbol-symmetric. Re-expressing the Beta functions in terms of Gamma functions, followed by some simplification, yields

$$\mathbb{E} \left[\frac{1}{P_Y} \right] = \frac{c}{1 - \kappa} \frac{\Gamma(2\kappa)}{\Gamma(\kappa)} \frac{\Gamma(c + \kappa - 1)}{\Gamma(c + 2\kappa - 1)} + \frac{1 - 2\kappa}{1 - \kappa}. \quad (80)$$

Due to the symbol symmetry of $f_{\mathbf{P}}$, the derivations for the all-1 attack are the same. ■

B. Saddle point analysis for the Interleaving score function

In this section we will show that the best attack against the interleaving defense is the interleaving attack. In particular, this means that using the interleaving defense attains capacity against any attack. We do this in two steps: first we will show the existence of a saddlepoint in the attack vs. distribution space, and then we will argue that for a fixed distribution this saddle point attains the minimum value of the performance indicator - in other words, that there is no better attack.

In the first step we do a saddlepoint analysis of the performance indicator $\tilde{\mu}_C / \tilde{\sigma}_{\text{inn}}$, in the following setting. We fix the employed suspicion function h to be the ‘Interleaving defense’ as specified in Proposition 9. The tracer has to tune the bias distribution $f_{\mathbf{P}}$ and at the same time the coalition has to find the best possible attack $f_{y|m}$ against the combination $f_{\mathbf{P}}, h$. This simultaneous counter-acting optimization leads to a saddle point solution for $\tilde{\mu}_C / \tilde{\sigma}_{\text{inn}}$ which is a minimum as a function of the attack strategy and a maximum as a function of $f_{\mathbf{P}}$. A similar analysis was done by Huang and Moulin [20] in the context of the q -ary fingerprinting capacity, abstracting away the exact suspicion function to be employed. They found the saddlepoint at (attack = interleaving, $f_{\mathbf{P}}$ = Dirichlet with $\kappa = \frac{1}{2}$), consistent with the asymptotic ($c \rightarrow \infty$) fingerprinting capacity $(q - 1) / (2c^2 \ln q)$ known earlier [19].

We use the Lagrangian approach, with functional L given by

$$L = \frac{\tilde{\mu}_C^2}{\tilde{\sigma}_{\text{inn}}^2} + \sum_{\mathbf{m}} \lambda_{\mathbf{m}} \left(\sum_{y \in \mathcal{Q}} f_{y|\mathbf{m}} - 1 \right) + \Lambda \left(\int d^q p \delta(|\mathbf{p}| - 1) f_{\mathbf{P}} - 1 \right). \quad (81)$$

Here the $\lambda_{\mathbf{m}}$ and Λ are constraint multipliers: $\lambda_{\mathbf{m}}$ multiplies the constraint that, for every \mathbf{m} , $f_{y|\mathbf{m}}$ is a probability mass function for y , and Λ multiplies the constraint that $f_{\mathbf{P}}$ is a probability density function.

Proposition 17. *Let the tracer use the interleaving defense. When the interleaving attack strategy is used, and the bias distribution is the Dirichlet distribution with $\kappa = \frac{1}{2}$, a saddlepoint occurs. Also, the asymptotic fingerprinting rate in the saddlepoint is equal to the asymptotic fingerprinting capacity $\frac{q-1}{2c^2 \ln q}$.*

Proof: The proof consists of two parts: (i) showing that we have a stationary point in the $(f_{y|m}, f_p)$ -space which corresponds to capacity; (ii) showing that the stationary point is a maximum as a function of f_p and a minimum as a function of $f_{y|m}$.

Part 1. For the interleaving defense we have from Lemma 11 and Corollary 6 that

$$\tilde{\mu}_C \sqrt{q-1} = \sum_{y \in \mathcal{Q}} \mathbb{E}_{\mathbf{P}} \left[\left. \frac{\partial f_{y|\mathbf{P}}}{\partial P_y} \right|_{|\mathbf{p}|=1} \right] \quad (82)$$

$$= -c + \sum_{y \in \mathcal{Q}} \sum_{\mathbf{m}} \binom{c}{\mathbf{m}} f_{y|\mathbf{m}} m_y \mathbb{E}_{\mathbf{P}} \left[\frac{\mathbf{P}^{\mathbf{m}}}{P_y} \right], \quad (83)$$

and

$$\tilde{\sigma}_{\text{inn}}^2 (q-1) = -1 + \sum_{y \in \mathcal{Q}} \mathbb{E}_{\mathbf{P}} \left[\frac{f_{y|\mathbf{P}}}{P_y} \right] \quad (84)$$

$$= -1 + \sum_{y \in \mathcal{Q}} \sum_{\mathbf{m}} \binom{c}{\mathbf{m}} f_{y|\mathbf{m}} \mathbb{E}_{\mathbf{P}} \left[\frac{\mathbf{P}^{\mathbf{m}}}{P_y} \right]. \quad (85)$$

The functional derivatives of $\tilde{\mu}_C$ and $\tilde{\sigma}_{\text{inn}}^2$ are

$$\sqrt{q-1} \frac{\delta \tilde{\mu}_C}{\delta f_{y|\mathbf{m}}} = \binom{c}{\mathbf{m}} m_y \mathbb{E}_{\mathbf{P}} \left[\frac{\mathbf{P}^{\mathbf{m}}}{P_y} \right]; \quad (86)$$

$$\sqrt{q-1} \frac{\delta \tilde{\mu}_C}{\delta f_p} = \delta(|\mathbf{p}|-1) \sum_{y \in \mathcal{Q}} \left. \frac{\partial f_{y|\mathbf{p}}}{\partial p_y} \right|_{|\mathbf{p}|=1}; \quad (87)$$

$$(q-1) \frac{\delta \tilde{\sigma}_{\text{inn}}^2}{\delta f_{y|\mathbf{m}}} = \binom{c}{\mathbf{m}} \mathbb{E}_{\mathbf{P}} \left[\frac{\mathbf{P}^{\mathbf{m}}}{P_y} \right]; \quad (88)$$

$$(q-1) \frac{\delta \tilde{\sigma}_{\text{inn}}^2}{\delta f_p} = \delta(|\mathbf{p}|-1) \sum_{y \in \mathcal{Q}} \frac{f_{y|\mathbf{p}}}{p_y}. \quad (89)$$

With these ingredients, the stationarity equations become

$$0 = \frac{\delta L}{\delta f_{y|\mathbf{m}}} = \lambda_{\mathbf{m}} + \frac{2\tilde{\mu}_C}{\tilde{\sigma}_{\text{inn}}^2} \frac{\delta \tilde{\mu}_C}{\delta f_{y|\mathbf{m}}} - \frac{\tilde{\mu}_C^2}{\tilde{\sigma}_{\text{inn}}^4} \frac{\delta \tilde{\sigma}_{\text{inn}}^2}{\delta f_{y|\mathbf{m}}} \quad (90)$$

$$= \lambda_{\mathbf{m}} + \frac{2\tilde{\mu}_C}{\tilde{\sigma}_{\text{inn}}^2} \binom{c}{\mathbf{m}} \mathbb{E}_{\mathbf{P}} \left[\frac{\mathbf{P}^{\mathbf{m}}}{P_y} \right] \left[\frac{m_y}{\sqrt{q-1}} - \frac{\tilde{\mu}_C}{2\tilde{\sigma}_{\text{inn}}^2 (q-1)} \right]; \quad (91)$$

$$0 = \frac{\delta L}{\delta f_p} = \Lambda + \frac{2\tilde{\mu}_C}{\tilde{\sigma}_{\text{inn}}^2} \frac{\delta \tilde{\mu}_C}{\delta f_p} - \frac{\tilde{\mu}_C^2}{\tilde{\sigma}_{\text{inn}}^4} \frac{\delta \tilde{\sigma}_{\text{inn}}^2}{\delta f_p} \quad (92)$$

$$= \Lambda + \frac{2\tilde{\mu}_C}{\tilde{\sigma}_{\text{inn}}^2 \sqrt{q-1}} \sum_{y \in \mathcal{Q}} \left. \frac{\partial f_{y|\mathbf{p}}}{\partial p_y} \right|_{|\mathbf{p}|=1} - \frac{\tilde{\mu}_C^2}{\tilde{\sigma}_{\text{inn}}^4 (q-1)} \sum_{y \in \mathcal{Q}} \frac{f_{y|\mathbf{p}}}{p_y}. \quad (93)$$

Equation (91) has to hold for all symbols y . This means that the expression $\mathbb{E}_{\mathbf{P}} \left[\frac{\mathbf{P}^{\mathbf{m}}}{P_y} \right] \left[\frac{m_y}{\sqrt{q-1}} - \frac{\tilde{\mu}_C}{2\tilde{\sigma}_{\text{inn}}^2 (q-1)} \right]$ has to be independent of y . This is a very complicated requirement on f_p and the attack; in general there is no easy way of solving

it. However, if we take the Dirichlet distribution for f_p , with parameter κ , then

$$\mathbb{E}_{\mathbf{P}} \left[\frac{\mathbf{P}^{\mathbf{m}}}{P_y} \right] = \frac{B(\kappa \mathbf{1}_q + \mathbf{m})}{B(\kappa \mathbf{1}_q)} \frac{\kappa q + c - 1}{m_y - (1 - \kappa)} \quad (94)$$

and it becomes possible to satisfy the independence requirement by demanding

$$\tilde{\mu}_C = 2\tilde{\sigma}_{\text{inn}}^2 (1 - \kappa) \sqrt{q-1}. \quad (95)$$

With this special relation between $\tilde{\mu}_C$ and $\tilde{\sigma}_{\text{inn}}^2$, (93) becomes

$$\forall \mathbf{p}: \quad 0 = \Lambda + 4(1 - \kappa) \sum_{y \in \mathcal{Q}} \left. \frac{\partial f_{y|\mathbf{p}}}{\partial p_y} \right|_{|\mathbf{p}|=1} - 4(1 - \kappa)^2 \sum_{y \in \mathcal{Q}} \frac{f_{y|\mathbf{p}}}{p_y}. \quad (96)$$

We have to find strategy parameters $f_{y|\mathbf{m}}$ that give rise to a function $f_{y|\mathbf{p}}$ that satisfies (96). We happen to know from (43) that the interleaving attack satisfies

$$\sum_{y \in \mathcal{Q}} \left. \frac{\partial f_{y|\mathbf{p}}}{\partial p_y} \right|_{|\mathbf{p}|=1} = q-1 \quad \text{and} \quad \sum_{y \in \mathcal{Q}} \frac{f_{y|\mathbf{p}}}{p_y} = q. \quad (97)$$

Thus (96) is satisfied if we take the interleaving attack and $\Lambda = 4(1 - \kappa)^2 q - 4(1 - \kappa)(q-1)$.

Next, we know that $\tilde{\sigma}_{\text{inn}}^2 = 1$ in case of a match since our suspicion function is normalized, and Proposition 9 tells us that $\tilde{\mu}_C = \sqrt{q-1}$ for the interleaving match. The relationship (95) can only hold if $\kappa = 1/2$.

Thus, we have a stationary point in which the interleaving attack is used and

$$\begin{cases} \tilde{\mu}_C = \sqrt{q-1}; \\ \tilde{\sigma}_{\text{inn}}^2 = 1; \\ \kappa = \frac{1}{2}; \\ \Lambda = -(q-2); \\ \lambda_{\mathbf{m}} = -\frac{1}{2} \binom{c}{\mathbf{m}} \frac{B(\kappa \mathbf{1}_q + \mathbf{m})}{B(\kappa \mathbf{1}_q)} \left(\frac{q}{2} + c - 1 \right). \end{cases} \quad (98)$$

We have a match with $\tilde{\mu}_C^2 / \tilde{\sigma}_{\text{inn}}^2 = q-1$, which corresponds to asymptotic capacity as described in Proposition 10.

Part 2. We take an arbitrary point in $(f_{y|m}, f_p)$ -space and consider the infinitesimal steps

$$\begin{cases} \hat{f}_{y|\mathbf{m}} = f_{y|\mathbf{m}} + \Delta_{y\mathbf{m}} \\ \hat{f}_p = f_p + \beta(\mathbf{p}) \end{cases} \quad (99)$$

with $\sum_y \Delta_{y\mathbf{m}} = 0$ and $\int d^q p \delta(|\mathbf{p}|-1) \beta(\mathbf{p}) = 0$. In the new point we write

$$\hat{\mu}_C = \tilde{\mu}_{(0)} + \tilde{\mu}_{(1)} + \tilde{\mu}_{(2)} \quad ; \quad \hat{\sigma}_{\text{inn}}^2 = \tilde{\sigma}_{(0)}^2 + \tilde{\sigma}_{(1)}^2 + \tilde{\sigma}_{(2)}^2 \quad (100)$$

where $\tilde{\mu}_{(0)} = \tilde{\mu}_C$ and $\tilde{\sigma}_{(0)}^2 = \tilde{\sigma}_{\text{inn}}^2$ refer to values in the original

point $(f_{y|m}, f_p)$, and we have defined

$$\begin{aligned} \tilde{\mu}_{(1)}\sqrt{q-1} &= \sum_{y \in \mathcal{Q}} \sum_{\mathbf{m}} \int d^q p \delta(|\mathbf{p}|-1) \binom{c}{\mathbf{m}} \frac{p^m}{p_y} \\ &\quad \cdot m_y [f_p \Delta_{ym} + \beta(\mathbf{p}) f_{y|m}]; \end{aligned} \quad (101)$$

$$\begin{aligned} \tilde{\sigma}_{(1)}^2(q-1) &= \sum_{y \in \mathcal{Q}} \sum_{\mathbf{m}} \int d^q p \delta(|\mathbf{p}|-1) \binom{c}{\mathbf{m}} \frac{p^m}{p_y} \\ &\quad \cdot [f_p \Delta_{ym} + \beta(\mathbf{p}) f_{y|m}]; \end{aligned} \quad (102)$$

$$\begin{aligned} \tilde{\mu}_{(2)}\sqrt{q-1} &= \sum_{y \in \mathcal{Q}} \sum_{\mathbf{m}} \int d^q p \delta(|\mathbf{p}|-1) \binom{c}{\mathbf{m}} \frac{p^m}{p_y} \\ &\quad \cdot m_y \beta(\mathbf{p}) \Delta_{ym}; \end{aligned} \quad (103)$$

$$\begin{aligned} \tilde{\sigma}_{(2)}^2(q-1) &= \sum_{y \in \mathcal{Q}} \sum_{\mathbf{m}} \int d^q p \delta(|\mathbf{p}|-1) \binom{c}{\mathbf{m}} \frac{p^m}{p_y} \\ &\quad \cdot \beta(\mathbf{p}) \Delta_{ym}. \end{aligned} \quad (104)$$

The subscript indicates the order of the small step. The maximum order is 2, since the expressions for $\tilde{\mu}_C$ and $\tilde{\sigma}_{\text{inn}}^2$ are linear in both f_p and $f_{y|m}$. We investigate the fraction $\tilde{\mu}_C^2/\tilde{\sigma}_{\text{inn}}^2$.

$$\frac{\tilde{\mu}_C^2}{\tilde{\sigma}_{\text{inn}}^2} = \frac{[\tilde{\mu}_{(0)} + \tilde{\mu}_{(1)} + \tilde{\mu}_{(2)}]^2}{\tilde{\sigma}_{(0)}^2 + \tilde{\sigma}_{(1)}^2 + \tilde{\sigma}_{(2)}^2} \quad (105)$$

$$= \frac{\tilde{\mu}_{(0)}^2}{\tilde{\sigma}_{(0)}^2} \cdot \frac{1 + 2\frac{\tilde{\mu}_{(1)}}{\tilde{\mu}_{(0)}} + 2\frac{\tilde{\mu}_{(2)}}{\tilde{\mu}_{(0)}} + \frac{[\tilde{\mu}_{(1)}]^2}{\tilde{\mu}_{(0)}^2} + \dots}{1 + \frac{\tilde{\sigma}_{(1)}^2}{\tilde{\sigma}_{(0)}^2} + \frac{\tilde{\sigma}_{(2)}^2}{\tilde{\sigma}_{(0)}^2}} \quad (106)$$

$$= \frac{\tilde{\mu}_{(0)}^2}{\tilde{\sigma}_{(0)}^2} \cdot \left[1 + 2\frac{\tilde{\mu}_{(1)}}{\tilde{\mu}_{(0)}} + 2\frac{\tilde{\mu}_{(2)}}{\tilde{\mu}_{(0)}} + \frac{[\tilde{\mu}_{(1)}]^2}{[\tilde{\mu}_{(0)}]^2} + \dots \right] \cdot \left[1 - \frac{\tilde{\sigma}_{(1)}^2}{\tilde{\sigma}_{(0)}^2} - \frac{\tilde{\sigma}_{(2)}^2}{\tilde{\sigma}_{(0)}^2} + \frac{[\tilde{\sigma}_{(1)}^2]^2}{[\tilde{\sigma}_{(0)}^2]^2} + \dots \right] \quad (107)$$

where the dots stand for higher order terms. In the last line we did a Taylor expansion of the denominator. By collecting equal order terms in (107) we obtain the first and second order components

$$\left[\frac{\tilde{\mu}_C^2}{\tilde{\sigma}_{\text{inn}}^2} \right]_{(1)} = \frac{\tilde{\mu}_{(0)}^2}{\tilde{\sigma}_{(0)}^2} \cdot \left[2\frac{\tilde{\mu}_{(1)}}{\tilde{\mu}_{(0)}} - \frac{\tilde{\sigma}_{(1)}^2}{\tilde{\sigma}_{(0)}^2} \right] \quad (108)$$

and

$$\begin{aligned} \left[\frac{\tilde{\mu}_C^2}{\tilde{\sigma}_{\text{inn}}^2} \right]_{(2)} &= \frac{\tilde{\mu}_{(0)}^2}{\tilde{\sigma}_{(0)}^2} \left[2\frac{\tilde{\mu}_{(2)}}{\tilde{\mu}_{(0)}} + \frac{[\tilde{\mu}_{(1)}]^2}{[\tilde{\mu}_{(0)}]^2} - \frac{\tilde{\sigma}_{(2)}^2}{\tilde{\sigma}_{(0)}^2} + \right. \\ &\quad \left. + \frac{[\tilde{\sigma}_{(1)}^2]^2}{[\tilde{\sigma}_{(0)}^2]^2} - 2\frac{\tilde{\mu}_{(1)}}{\tilde{\mu}_{(0)}} \frac{\tilde{\sigma}_{(1)}^2}{\tilde{\sigma}_{(0)}^2} \right] \\ &= \frac{\tilde{\mu}_{(0)}^2}{\tilde{\sigma}_{(0)}^2} \left[2\frac{\tilde{\mu}_{(2)}}{\tilde{\mu}_{(0)}} - \frac{\tilde{\sigma}_{(2)}^2}{\tilde{\sigma}_{(0)}^2} \right] + \frac{\tilde{\mu}_{(0)}^2}{\tilde{\sigma}_{(0)}^2} \left[\frac{\tilde{\mu}_{(1)}}{\tilde{\mu}_{(0)}} - \frac{\tilde{\sigma}_{(1)}^2}{\tilde{\sigma}_{(0)}^2} \right]^2. \end{aligned} \quad (109)$$

We take the stationary point as our starting point and first make a step in the f_p -direction only, i.e. $\Delta_{ym} = 0$.

The fact that $\Delta_{ym} = 0$ yields $\tilde{\mu}_{(2)} = 0$ and $\tilde{\sigma}_{(2)}^2 = 0$ from (103) and (104). Furthermore, in (101) and (102) note that the sum over y yields a constant as in (97), and then integrating β gives zero. So $\tilde{\mu}_{(1)} = 0$ and $\tilde{\sigma}_{(1)}^2 = 0$ as well. Thus we conclude that from the stationary point, changing only f_p does not change the performance indicator $\tilde{\mu}_C^2/\tilde{\sigma}_{\text{inn}}^2$. Note that this is consistent with Proposition 10.

Secondly we fix f_p to be the Dirichlet distribution with $\kappa = \frac{1}{2}$ and vary the attack slightly from interleaving. Now $\beta = 0$, which yields $\tilde{\mu}_{(2)} = 0$ and $\tilde{\sigma}_{(2)}^2 = 0$. Equation (110) with $\beta = 0$ then reduces to a square, which is non-negative. Thus the performance indicator is minimized when the interleaving attack is used, and the found stationary point is indeed a saddlepoint. ■

This saddlepoint leads to a global minimum:

Theorem 2. *Assume the tracer uses the interleaving defense and the Dirichlet distribution with $\kappa = \frac{1}{2}$. Then the interleaving attack minimizes the performance indicator $\frac{\tilde{\mu}_C}{\tilde{\sigma}_{\text{inn}}}$.*

Proof: From Proposition 17 we know that the interleaving attack is a local minimum in this setting. Also, when the distribution is fixed as the Dirichlet distribution with $\kappa = \frac{1}{2}$, the proof of Proposition 17 states that the second derivative (110) is non-negative for any strategy, as $\beta(\mathbf{p}) = 0$ implies that $\tilde{\mu}_{(2)} = \tilde{\sigma}_{(2)} = 0$. Since $\frac{\tilde{\mu}_C}{\tilde{\sigma}_{\text{inn}}}$ is a rational function of Δ_{ym} , we can conclude that the interleaving attack is a *global* minimum for this setting. ■

C. Relation to the Tardos suspicion function

The interleaving defense is closely related to the Tardos suspicion function:

Proposition 18. *The symmetric Tardos function is the strongly normalized optimal suspicion function against the interleaving attack.*

Proof: We know from (48) that $\text{Var}_{X|p}[T] = \frac{1-p_y}{c^2 p_y}$. So by Theorem 8, the strongly normalized optimal suspicion function against the interleaving attack is

$$h(x, y, \mathbf{p}) = \sqrt{\frac{p_y}{1-p_y}} \left(\frac{\delta_{x,y}}{p_y} - 1 \right), \quad (111)$$

which equals the the symmetric Tardos function (3). ■

D. Interleaving Defense Numerics

To verify our analytic results and their practical applicability, we ran simulations for the binary case and the arcsine distribution (without cut-off), which is equal to the Dirichlet distribution with $\kappa = \frac{1}{2}$. We simulated the five described attacks (interleaving, all-1, coin-flip, majority voting, and minority voting) against the interleaving defense. We stress that these five attacks are by no means exhaustive.

We ran simulations for $1 \leq c \leq 200$ to obtain the $\tilde{\mu}_c$ and the $\tilde{\sigma}_{\text{inn}}$ in these five cases as depicted in Fig.2. We then analyzed this data to obtain the leading-order term in c . The results can be found in Table I. Since for mismatches the innocent score is no longer normalized ($\tilde{\sigma}_{\text{inn}} \neq 1$), we present the results for $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}}$ to make a fair comparison.

As predicted by Theorem 2, the interleaving defense attains capacity ($\tilde{\mu} = 1$) against the interleaving attack. We also observe that the majority voting attack has a constant $\tilde{\mu}$. For the other three attacks, $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}}$ seems to grow as $c^{1/4}$. We were able to prove this for the all-1 and coin-flip attacks in Proposition 16.

TABLE I

Numerical trends for the performance indicator $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}}$ of the interleaving defense in the binary case $q = 2$ for large c .

interleaving attack	1.0
all-1 attack	$0.61c^{0.23}$
coin-flip attack	$0.61c^{0.23}$
majority voting attack	1.2
minority voting attack	$0.75c^{0.25}$

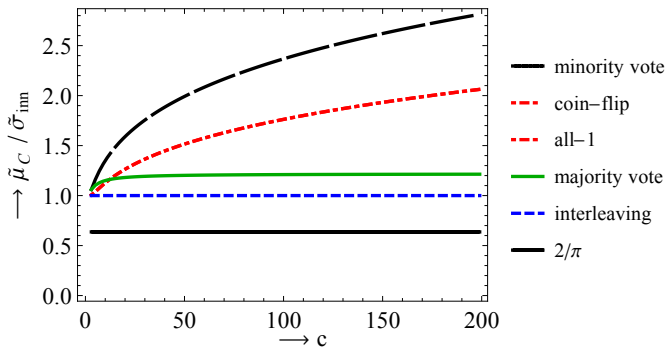


Fig. 2. Interleaving defense against various attacks in the binary case.

VI. ALL-HIGH DEFENSE

A. Optimal defense

The all-high attack

$$f_{y|m} = \begin{cases} 1 & \text{if } m_y > 0 \text{ and } m_{y+1} = \dots = m_{q-1} = 0 \\ 0 & \text{else} \end{cases} \quad (112)$$

outputs the highest symbol among those received by the coalition.

Note that this is the only attack we consider that breaks symbol symmetry and assumes an ordering of the alphabet. This is a special case of the so-called *preferred-sequence attack*, in which the colluders have a predetermined ranking of the symbols. The results below generalize to the preferred-sequence attack. Recall our shorthand notation $a_k := (p_0 + \dots + p_{k-1})$ and $a_{\mathcal{B}} = \sum_{\beta \in \mathcal{B}} p_{\beta}$.

Proposition 19. *Against the all-high attack, the optimal suspicion function is $h = (T - 1)/\sqrt{\text{Var}[T]}$, with*

$$T(x, y, \mathbf{p}) = \begin{cases} (a_{y+1}^{c-1} - a_y^{c-1})/(a_{y+1}^c - a_y^c) & \text{if } x < y \\ a_{y+1}^{c-1}/(a_{y+1}^c - a_y^c) & \text{if } x = y \\ 0 & \text{if } x > y. \end{cases} \quad (113)$$

In case of a match, it holds that

$$\tilde{\mu}_c = c \sqrt{-1 + \mathbb{E}_{\mathbf{P}} \left[\sum_{y=0}^{q-1} \frac{A_{y+1}^{2c-1} - 2A_y^c A_{y+1}^{c-1} + A_y^{2c-1}}{A_{y+1}^c - A_y^c} \right]}. \quad (114)$$

Proof: We find

$$f_{y|\mathbf{p}} = \mathbb{E}_{M|\mathbf{p}}[f_{y|M}] \quad (115)$$

$$= \mathbb{P}[M_y > 0, M_{y+1} = \dots = M_{q-1} = 0] \quad (116)$$

$$= \mathbb{P}[M_{y+1} = \dots = M_{q-1} = 0] \quad (117)$$

$$- \mathbb{P}[M_y = \dots = M_{q-1} = 0]$$

$$= \frac{a_{y+1}^c}{|\mathbf{p}|^c} - \frac{a_y^c}{|\mathbf{p}|^c}. \quad (118)$$

So

$$T = \frac{1}{c} \frac{\partial \ln(|\mathbf{p}|^c f_{y|\mathbf{p}})}{\partial p_x} \Big|_{|\mathbf{p}|=1} \quad (119)$$

$$= \begin{cases} \frac{a_{y+1}^{c-1} - a_y^{c-1}}{a_{y+1}^c - a_y^c} & \text{if } x < y \\ \frac{a_{y+1}^{c-1}}{a_{y+1}^c - a_y^c} & \text{if } x = y \\ 0 & \text{if } x > y. \end{cases} \quad (120)$$

Also,

$$\mathbb{E}[T^2] = \mathbb{E}_{\mathbf{P}} \mathbb{E}_{Y|\mathbf{P}} \mathbb{E}_{X|\mathbf{P}} [T^2(X, Y, \mathbf{P})] \quad (121)$$

$$= \mathbb{E}_{\mathbf{P}} \mathbb{E}_{Y|\mathbf{P}} [P_Y T^2(Y, Y, \mathbf{P}) + A_Y T^2(0, Y, \mathbf{P})] \quad (122)$$

$$= \mathbb{E}_{\mathbf{P}} \sum_{y=0}^{q-1} \left[P_y \frac{A_{y+1}^{2(c-1)}}{A_{y+1}^c - A_y^c} + A_y \frac{(A_{y+1}^{c-1} - A_y^{c-1})^2}{A_{y+1}^c - A_y^c} \right] \quad (123)$$

$$= \mathbb{E}_{\mathbf{P}} \sum_{y=0}^{q-1} \frac{A_{y+1}^{2c-1} - 2A_y^c A_{y+1}^{c-1} + A_y^{2c-1}}{A_{y+1}^c - A_y^c}. \quad (124)$$

We obtain (114) using $\tilde{\mu}_c = c\sqrt{\text{Var}[T]} = c\sqrt{\mathbb{E}[T^2] - 1}$. ■

When $x = y$, the h is positive. When $x > y$, it is negative and constant. When $x < y$, it might be negative or it might not. For instance, for $c = 2$, we find $(a_{y+1} - a_y)/(a_{y+1}^2 - a_y^2) = 1/(a_{y+1} + a_y) = 1/(p_y + 2a_y)$, in which case h is negative if and only if $p_y > 1 - 2a_y$. In particular it is negative if $a_y \geq \frac{1}{2}$. Also, h is the same for all $x < y$.

We now analyze the behaviour of $\tilde{\mu}_c$ when the symmetric Dirichlet distribution is employed. Before we can state our result, we will need the following Lemma:

Lemma 20. *Let $f_{\mathbf{P}}$ be the symmetric Dirichlet distribution without cutoff. The joint distribution for the pair*

$(A_{y+1}, A_y/A_{y+1})$ is then given by

$$J(a_{y+1}, \frac{a_y}{a_{y+1}}) = \frac{a_{y+1}^{-1+(y+1)\kappa} (1 - a_{y+1})^{-1+(q-y-1)\kappa}}{B([y+1]\kappa, [q-y-1]\kappa)} \times \frac{(a_y/a_{y+1})^{-1+y\kappa} (1 - a_y/a_{y+1})^{-1+\kappa}}{B(y\kappa, \kappa)}.$$

Proof: We first derive the joint distribution $J(a_y, a_{y+1})$ for A_y and A_{y+1} :

$$J(a_y, a_{y+1}) = \mathbb{E}_{\mathcal{P}} \left[\delta \left[A_y - \sum_{i=0}^{y-1} P_i \right] \delta \left[A_{y+1} - \sum_{i=0}^y P_i \right] \right] \quad (125)$$

$$\propto \int_{|\mathbf{p}|=1} d^{q-1} \mathbf{p} \mathbf{p}^{\kappa-1} \delta \left[a_y - \sum_{i=0}^{y-1} p_i \right] \delta \left[a_{y+1} - \sum_{i=0}^y p_i \right] \quad (126)$$

$$= \int d^q \mathbf{p} \mathbf{p}^{\kappa-1} \delta \left[a_y - \sum_{i=0}^{y-1} p_i \right] \delta \left[a_{y+1} - \sum_{i=0}^y p_i \right] \delta(1 - |\mathbf{p}|) \quad (127)$$

$$= \int d^q \mathbf{p} \mathbf{p}^{\kappa-1} \delta \left[a_y - \sum_{i=0}^{y-1} p_i \right] \delta \left[a_{y+1} + \sum_{i=y+1}^{q-1} p_i - 1 \right] \delta(1 - |\mathbf{p}|). \quad (128)$$

Here $\delta(x)$ is the Dirac delta function. We perform the following change of variables: for $i < y$ we define $p_i = a_y s_i$; for $i > y$ we define $p_i = (1 - a_{y+1}) t_i$. This yields $d^q \mathbf{p} \mathbf{p}^{\kappa-1} = d^y \mathbf{s} d^y \mathbf{t} d^{q-y-1} t p_y^{\kappa-1} a_y^{y\kappa} \mathbf{s}^{\kappa-1} (1 - a_{y+1})^{[q-y-1]\kappa} \mathbf{t}^{\kappa-1}$ and

$$\delta(a_y - \sum_{i=0}^{y-1} p_i) = a_y^{-1} \delta(1 - |\mathbf{s}|), \quad (129)$$

$$\delta \left[a_{y+1} + \sum_{i=y+1}^{q-1} p_i - 1 \right] = (1 - a_{y+1})^{-1} \delta(1 - |\mathbf{t}|), \quad (130)$$

$$\delta(1 - |\mathbf{p}|) = \delta[1 - p_y - a_y |\mathbf{s}| - (1 - a_{y+1}) |\mathbf{t}|]. \quad (131)$$

The expression (128) becomes

$$J(a_y, a_{y+1}) = \int d^y \mathbf{s} d^y \mathbf{t} d^{q-y-1} t p_y^{\kappa-1} a_y^{y\kappa-1} \mathbf{s}^{\kappa-1} \times (1 - a_{y+1})^{[q-y-1]\kappa-1} \mathbf{t}^{\kappa-1} \delta(1 - |\mathbf{s}|) \delta(1 - |\mathbf{t}|) \delta[p_y + a_y - a_{y+1}] \quad (132) \\ \propto a_y^{y\kappa-1} (1 - a_{y+1})^{[q-y-1]\kappa-1} (a_{y+1} - a_y)^{\kappa-1}. \quad (133)$$

Finally we do a last change of variables from a_y to $z = a_y/a_{y+1}$. This gives $da_y da_{y+1} = a_{y+1} da_{y+1} dz$, and (133) becomes

$$J(a_{y+1}, z) \propto a_{y+1}^{[y+1]\kappa-1} (1 - a_{y+1})^{[q-y-1]\kappa-1} z^{y\kappa-1} (1 - z)^{\kappa-1}. \quad (134)$$

Inserting the normalization constants yields the result of the lemma. ■

Given this joint distribution, we can now derive our main result for the all-high attack when the symmetric Dirichlet distribution is used.

Proposition 21. *Let $f_{\mathcal{P}}$ be the symmetric Dirichlet distribution without cutoff. If the attack is the all-high attack and the defense matches it, then, for large c ,*

$$\tilde{\mu}_C = c^{\frac{1-\kappa}{2}} \sqrt{\frac{\kappa \Gamma(q\kappa) \zeta(1 + \kappa)}{\Gamma([q-1]\kappa)}} \left[1 + \mathcal{O}(c^{-\min(1, \kappa)}) \right], \quad (135)$$

where ζ is the Riemann zeta function.

Proof: We write (114) as

$$\frac{\tilde{\mu}_C^2}{c^2} = -1 + \sum_{y=0}^{q-1} \mathbb{E}_{\mathcal{P}} \left[A_{y+1}^{c-1} \frac{1 - 2(A_y/A_{y+1})^c + (A_y/A_{y+1})^{2c-1}}{1 - (A_y/A_{y+1})^c} \right]. \quad (136)$$

The fraction can be expanded as

$$\frac{1}{1 - (A_y/A_{y+1})^c} = \sum_{t=0}^{\infty} (A_y/A_{y+1})^{tc}. \quad (137)$$

Then we evaluate the expectation using the joint distribution $J(a_{y+1}, \frac{a_y}{a_{y+1}})$ from Lemma 20. This yields

$$\frac{\tilde{\mu}_C^2}{c^2} = -1 + \sum_{y=0}^{q-1} \frac{B([y+1]\kappa + c - 1, [q-y-1]\kappa)}{B([y+1]\kappa, [q-y-1]\kappa)} \times \left[1 - \sum_{t=1}^{\infty} \frac{B(y\kappa + tc, \kappa)}{B(y\kappa, \kappa)} + \sum_{t=2}^{\infty} \frac{B(y\kappa + tc - 1, \kappa)}{B(y\kappa, \kappa)} \right] \quad (138)$$

noting that $1/B(y\kappa, \kappa)$ vanishes for $y = 0$. Further simplification gives

$$\frac{\tilde{\mu}_C^2}{c^2} = \frac{\kappa \Gamma(q\kappa)}{\Gamma(q\kappa + c - 1)} \left[\sum_{y=0}^{q-2} \frac{\Gamma([y+1]\kappa + c - 1)}{([y+2]\kappa + c - 1) \Gamma([y+1]\kappa)} + \sum_{y=1}^{q-1} \frac{\Gamma([y+1]\kappa + c - 1)}{\Gamma(y\kappa)} \sum_{t=2}^{\infty} \frac{\Gamma(y\kappa + tc - 1)}{\Gamma([y+1]\kappa + tc)} \right]. \quad (139)$$

Finally we use the identity $\Gamma(c+a)/\Gamma(c+b) = c^{a-b} [1 + \mathcal{O}(c^{-1})]$ to investigate the asymptotics. In the first summation over y the dominant term occurs at $y = q-2$, thus the summation can be simplified to $c^{-\kappa-1} [1 + \mathcal{O}(c^{-\min(1, \kappa)})]$. $\kappa \Gamma(q\kappa)/\Gamma([q-1]\kappa)$. Similarly, in the second summation over y the dominant term occurs at $y = q-1$ and thus this summation reduces to $c^{-\kappa-1} [1 + \mathcal{O}(c^{-\min(1, \kappa)})] [\zeta(1 + \kappa) - 1] \kappa \Gamma(q\kappa)/\Gamma([q-1]\kappa)$, where ζ is the Riemann zeta function. ■

B. All-1 defense

The binary all-high attack is known as the all-1 attack. It has $f_{1|m} = 1$ whenever $m_1 > 0$ and $f_{1|m} = 0$ when $m_1 = 0$.

Corollary 22. *Against the all-1 attack, the optimal suspicion function is $h = (T - 1)/\sqrt{\text{Var}[T]}$, with*

$$T(x, y, \mathbf{p}) = \begin{cases} (1 - p_0^{c-1})/(1 - p_0^c) & \text{if } (x, y) = (0, 1) \\ 1/(1 - p_0^c) & \text{if } (x, y) = (1, 1) \\ 1/p_0 & \text{if } (x, y) = (0, 0) \\ 0 & \text{if } (x, y) = (1, 0). \end{cases} \quad (140)$$

In case of a match it holds that

$$\tilde{\mu}_C = c\sqrt{\mathbb{E}_{\mathbf{P}}[P_0^{c-1}(1 - P_0)/(1 - P_0^c)]}. \quad (141)$$

Proof: T follows directly from Proposition 19. Furthermore, (114) gives

$$\text{Var}[T] = -1 + \mathbb{E}_{\mathbf{P}} \left[P_0^{c-1} + \frac{1 - 2P_0^c + P_0^{2c-1}}{1 - P_0^c} \right] \quad (142)$$

$$= \mathbb{E}_{\mathbf{P}} \left[P_0^{c-1} + \frac{P_0^{2c-1} - P_0^c}{1 - P_0^c} \right] = \mathbb{E}_{\mathbf{P}} \left[\frac{P_0^c(1 - P_0)}{P_0(1 - P_0^c)} \right]. \quad (143)$$

When $x < y$, the h is positive for any c , in contrast to the q -ary case.

Corollary 23. *Let $f_{\mathbf{P}}$ be the symmetric Dirichlet distribution with $\kappa = \frac{1}{2}$ and cutoff $\delta = 0$. Against the all-1 attack, the optimal suspicion function attains $\tilde{\mu}_C \propto c^{1/4}$ for large c .*

Before we investigate the behaviour of the all-high defense against an interleaving attack, we first prove a general lemma about the interleaving attack.

Lemma 24. *If the tracer uses a strongly centered score function and the coalition uses the interleaving attack, then*

$$\tilde{\mu}_C = \sum_{y \in \mathcal{A}} \mathbb{E}_{\mathbf{P}}[P_y h(y, y, \mathbf{P})]. \quad (144)$$

Proof: For the interleaving attack, $cT = \frac{\delta_{x,y}}{p_y} + c - 1$, so

$$\tilde{\mu}_C = c\mathbb{E}[T \cdot h] \quad (145)$$

$$= \mathbb{E}_{\mathbf{P}} \mathbb{E}_{Y|\mathbf{P}} \mathbb{E}_{X|\mathbf{P}} \left[\left(\frac{\delta_{X,Y}}{P_Y} + c - 1 \right) h(X, Y, \mathbf{P}) \right] \quad (146)$$

$$= \mathbb{E}_{\mathbf{P}} \mathbb{E}_{Y|\mathbf{P}} \mathbb{E}_{X|\mathbf{P}} \left[\frac{\delta_{X,Y}}{P_Y} h(X, Y, \mathbf{P}) \right] \quad (147)$$

$$= \mathbb{E}_{\mathbf{P}} \mathbb{E}_{Y|\mathbf{P}} [h(Y, Y, \mathbf{P})]. \quad (148)$$

where (148) holds since $\mathbb{E}[h] = 0$. ■

The performance of the all-high defense against the interleaving attack can be analyzed as follows:

Proposition 25. *If the tracer uses the all-high defense but the coalition uses the interleaving attack, then*

$$\tilde{\mu}_C = \frac{1}{\sqrt{\text{Var}[T]}} \mathbb{E}_{\mathbf{P}} \left[\sum_{y=1}^{q-1} \frac{P_y A_{y+1}^{c-1}}{A_{y+1}^c - A_y^c} \right] \quad (149)$$

where T belongs to the all-high defense.

Proof: Applying Lemma 24 we obtain

$$\tilde{\mu}_C = \frac{1}{\sqrt{\text{Var}[T]}} \left(-1 + \mathbb{E}_{\mathbf{P}} \left[\sum_{y \in \mathcal{A}} \frac{P_y A_{y+1}^{c-1}}{A_{y+1}^c - A_y^c} \right] \right) \quad (150)$$

$$= \frac{1}{\sqrt{\text{Var}[T]}} \mathbb{E}_{\mathbf{P}} \left[\sum_{y=1}^{q-1} \frac{P_y A_{y+1}^{c-1}}{A_{y+1}^c - A_y^c} \right]. \quad (151)$$

■

In the binary case this reduces to

Proposition 26. *For $q = 2$, if the tracer uses the all-1 defense, but the coalition uses the interleaving attack, then*

$$\tilde{\mu}_C = \frac{1}{\sqrt{\text{Var}[T]}} \mathbb{E}_{\mathbf{P}} \left[P_1 \sum_{k=0}^{\infty} P_0^{kc} \right]. \quad (152)$$

Proof: Applying Lemma 24 we obtain

$$\tilde{\mu}_C = \frac{1}{\sqrt{\text{Var}[T]}} \left(-1 + \mathbb{E}_{\mathbf{P}} \left[1 + \frac{P_1}{1 - P_0^c} \right] \right). \quad (153)$$

■

The scaling behaviour for large c is

Lemma 27. *Let $q = 2$ and $f_{\mathbf{P}}$ be the symmetric Dirichlet distribution with parameter $\kappa = \frac{1}{2}$ without cutoff. If the tracer uses the all-1 defense, but the coalition uses the interleaving attack, then*

$$\tilde{\mu}_C = \frac{\Gamma(\kappa + 1)}{B(\kappa, \kappa) \sqrt{\text{Var}[T]}} \sum_{t=0}^{\infty} \frac{\Gamma(tc + \kappa)}{\Gamma(tc + 2\kappa + 1)}. \quad (154)$$

For large c , this scales as $c^{(\kappa+1)/2}$.

C. All-1 Defense Numerics

We ran simulations for the binary case and the arcsine distribution (without cut-off) with the same parameters as described in Section V-D. The table looks very similar to that of the interleaving defense. As expected, the all-1 defense performs better against the all-1 attack, but worse against the other four attacks. However, it retains the same scaling behaviour.

We again stress that these five attacks are by no means exhaustive.

TABLE II
Numerical trends for the performance indicator $\tilde{\mu}_C/\tilde{\sigma}_{\text{inn}}$ of the all-1 defense in the binary case $q = 2$ for large c .

interleaving attack	0.71
all-1 attack	$0.86c^{0.25}$
coin-flip attack	$0.44c^{0.23}$
majority voting attack	0.84
minority voting attack	$0.54c^{0.25}$

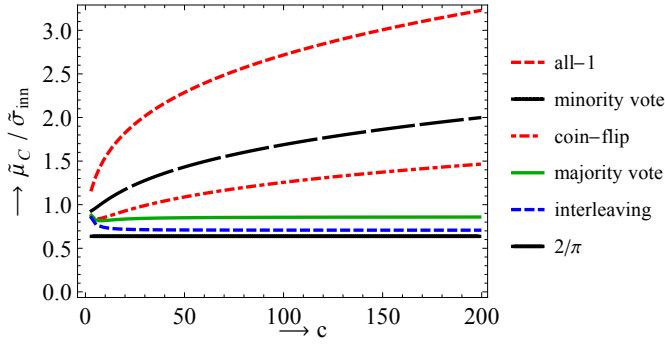


Fig. 3. All-1 defense against various attacks in the binary case.

VII. RANDOM-SYMBOL DEFENSE

A. Optimal defense

The random-symbol attack selects one of the received symbols uniformly at random. Tallies are disregarded, but a symbol can only be chosen if its tally is nonzero. The attack is parametrized by $f_{y|m} = (1 - \delta_{m_y,0})/|\{\alpha \in \mathcal{A} : m_\alpha > 0\}|$.

Proposition 28. For the random-symbol attack we find

$$|\mathbf{p}|^c f_{y|\mathbf{p}} = \frac{a_{\mathcal{A}}^c - a_{\mathcal{A} \setminus \{y\}}^c}{q} + \sum_{\mathcal{B} \subsetneq \mathcal{A}: y \in \mathcal{B}} \frac{a_{\mathcal{B}}^c - a_{\mathcal{B} \setminus \{y\}}^c}{|\mathcal{B}|(|\mathcal{B}| + 1)}. \quad (155)$$

The optimal suspicion function is $h = (T - 1)/\sqrt{\text{Var}[T]}$, with

$$T(x, y, \mathbf{p}) = \frac{1}{c} \frac{\partial \ln(|\mathbf{p}|^c f_{y|\mathbf{p}})}{\partial p_x} \Big|_{|\mathbf{p}|=1} = \quad (156)$$

$$\begin{cases} \frac{1}{f_{y|\mathbf{p}}} \left(\frac{1}{q} + \sum_{\mathcal{B} \subsetneq \mathcal{A}: y \in \mathcal{B}} \frac{a_{\mathcal{B}}^{c-1}}{|\mathcal{B}|(|\mathcal{B}| + 1)} \right) & \text{if } x = y \\ \frac{1}{f_{y|\mathbf{p}}} \left(\frac{1 - (1 - p_y)^{c-1}}{q} + \sum_{\substack{\mathcal{B} \subsetneq \mathcal{A} \\ x, y \in \mathcal{B}}} \frac{a_{\mathcal{B}}^{c-1} - a_{\mathcal{B} \setminus \{y\}}^{c-1}}{|\mathcal{B}|(|\mathcal{B}| + 1)} \right) & \text{if } x \neq y \end{cases} \quad (157)$$

Proof: For the random-symbol attack, the probability $f_{y|m}$ that the symbol y is produced, is 0 if $m_y = 0$. It is $\frac{1}{q}$ if for all $\alpha \in \mathcal{A}$, $m_\alpha > 0$. It is $\frac{1}{q-1}$ if $m_y > 0$ and there is exactly one symbol $\alpha_1 \in \mathcal{A}$ for which $m_{\alpha_1} = 0$. It is $\frac{1}{q-2}$ if $m_y > 0$ and there are exactly two distinct symbols $\alpha_1, \alpha_2 \in \mathcal{A}$ for which $m_{\alpha_1} = m_{\alpha_2} = 0$, etc. This can be written in additive form using indicator functions:

$$\begin{aligned} f_{y|m} &= \frac{1}{q} \mathbf{1}_{\{m_y > 0\}} \\ &+ \left(\frac{1}{q-1} - \frac{1}{q} \right) \mathbf{1}_{\{m_y > 0\}} \mathbf{1}_{\{\exists \alpha_1: m_{\alpha_1} = 0\}} \\ &+ \left(\frac{1}{q-2} - \frac{1}{q-1} \right) \mathbf{1}_{\{m_y > 0\}} \mathbf{1}_{\{\exists \alpha_1: m_{\alpha_1} = 0\}} \mathbf{1}_{\{\exists \alpha_2 \neq \alpha_1: m_{\alpha_2} = 0\}} \\ &+ \dots + \left(1 - \frac{1}{2} \right) \mathbf{1}_{\{m_y > 0\}} \\ &\cdot \mathbf{1}_{\{\exists \alpha_1: m_{\alpha_1} = 0\}} \dots \mathbf{1}_{\{\exists \alpha_{q-1} \neq \alpha_1, \dots, \alpha_{q-2}: m_{\alpha_{q-1}} = 0\}}. \end{aligned} \quad (158)$$

Note that

$$\mathbb{P}[M_y > 0] = \mathbb{P}[M_y \geq 0] - \mathbb{P}[M_y = 0] = \frac{A_{\mathcal{A}}^c - A_{\mathcal{A} \setminus \{y\}}^c}{|\mathbf{p}|^c} \quad (159)$$

and for each proper subset $\mathcal{B} \subsetneq \mathcal{A}$ with $y \in \mathcal{B}$, it holds that

$$\mathbb{P}[\forall \beta \in \mathcal{B}, M_\beta > 0] = \left(A_{\mathcal{B}}^c - A_{\mathcal{B} \setminus \{y\}}^c \right) / |\mathbf{p}|^c. \quad (160)$$

Since $f_{y|\mathbf{p}} = \mathbb{E}_{M|\mathbf{p}}[f_{y|M}]$, and for all sets \mathcal{V}, \mathcal{W} , it holds that $\mathbf{1}_{\mathcal{V}} \mathbf{1}_{\mathcal{W}} = \mathbf{1}_{\mathcal{V} \cap \mathcal{W}}$, and $\mathbb{E}[\mathbf{1}_{\mathcal{V}}] = \mathbb{P}[\mathcal{V}]$, we find

$$\begin{aligned} |\mathbf{p}|^c f_{y|\mathbf{p}} &= \frac{a_{\mathcal{A}}^c - a_{\mathcal{A} \setminus \{y\}}^c}{q} \\ &+ \sum_{\mathcal{B} \subsetneq \mathcal{A}: y \in \mathcal{B}} \left(\frac{1}{|\mathcal{B}|} - \frac{1}{|\mathcal{B}| + 1} \right) \left(a_{\mathcal{B}}^c - a_{\mathcal{B} \setminus \{y\}}^c \right). \end{aligned} \quad (161)$$

which simplifies to equation (155). \blacksquare

B. Coin-flip defense

The binary random-symbol attack is known as the coin-flip attack, and is parametrized as $f_{y|m} = \frac{1}{2}(1 - \delta_{m_y,0} + \delta_{m_y,c})$.

Proposition 29. Against the coin-flip attack, the optimal suspicion function is $h = (T - 1)/\sqrt{\text{Var}[T]}$, with

$$T(x, y, \mathbf{p}) = \begin{cases} (1 + p_y^{c-1}) / (1 + p_y^c - p_{1-y}^c) & \text{if } x = y \\ (1 - p_{1-y}^{c-1}) / (1 + p_y^c - p_{1-y}^c) & \text{if } x \neq y. \end{cases} \quad (162)$$

Proof: Since

$$f_{y|m} = \frac{1}{2}(1 - \delta_{m_y,0} + \delta_{m_y,c}). \quad (163)$$

$$f_{y|\mathbf{p}} = \mathbb{E}_{M|\mathbf{p}}[f_{y|M}] \quad (164)$$

$$= \frac{1}{2|\mathbf{p}|^c} \sum_{m_y=0}^c \binom{c}{m_y} p_y^{m_y} p_{1-y}^{c-m_y} (1 - \delta_{m_y,0} + \delta_{m_y,c}) \quad (165)$$

$$= \frac{1}{2|\mathbf{p}|^c} [(p_y + p_{1-y})^c - p_{1-y}^c + p_y^c]. \quad (166)$$

Thus

$$\begin{aligned} \frac{\partial(|\mathbf{p}|^c f_{y|\mathbf{p}})}{\partial p_x} &= \frac{1}{2} c [(p_y + p_{1-y})^{c-1} - (1 - \delta_{x,y}) p_{1-y}^{c-1} \\ &+ \delta_{x,y} p_y^{c-1}] \end{aligned} \quad (167)$$

So

$$T = \frac{1 - (1 - \delta_{x,y}) p_{1-y}^{c-1} + \delta_{x,y} p_y^{c-1}}{1 - p_{1-y}^c + p_y^c}. \quad (168)$$

When $x = y$, the h is positive. When $x \neq y$, it is negative, since $-p_{1-y}^{c-1} < p_y^{c-1}$, so $p_{1-y}^{c-1}(p_{1-y} - 1) < p_y^c$, and thus $1 - p_{1-y}^{c-1} < 1 + p_y^c - p_{1-y}^c$.

The interleaving attack against the coin-flip defense behaves as follows in the binary case:

Lemma 30. For $q = 2$, if the tracer uses the coin-flip defense, but the coalition uses the interleaving attack, then

$$\tilde{\mu}_c = \frac{1}{\sqrt{\text{Var}[T]}} \left[-1 + \mathbb{E}_{\mathcal{P}} \left[\frac{P_0(1+P_0)^{c-1}}{1+P_0^c - P_1^c} + \frac{P_1(1+P_1)^{c-1}}{1+P_1^c - P_0^c} \right] \right]. \quad (169)$$

Proof: This follows directly from Lemma 24 with (162). ■

C. Coin-flip Defense Numerics

We ran simulations for the binary case and the arcsine distribution (without cut-off) with the same parameters as described in Section V-D. The results look quite different to those of the interleaving defense. As expected, the coin-flip defense performs better against the coin-flip attack. However, against a majority voting attack this defense fails, as no information on the coalition is gained. There is still a small advantage left against an interleaving attack. However, it retains the same scaling behaviour against the minority voting and all-1 attacks.

We note that the all-1 and coin-flip attacks numerically perform the same against this defense. We could only prove this fact analytically for the interleaving defense.

We again stress that these five attacks are by no means exhaustive.

TABLE III

Numerical trends for the performance indicator $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}}$ of the coin-flip defense in the binary case $q = 2$ for large c .

interleaving attack	$5.1c^{-0.71}$
all-1 attack	$0.72c^{0.25}$
coin-flip attack	$0.72c^{0.25}$
majority voting attack	0.0
minority voting attack	$1.1c^{0.25}$

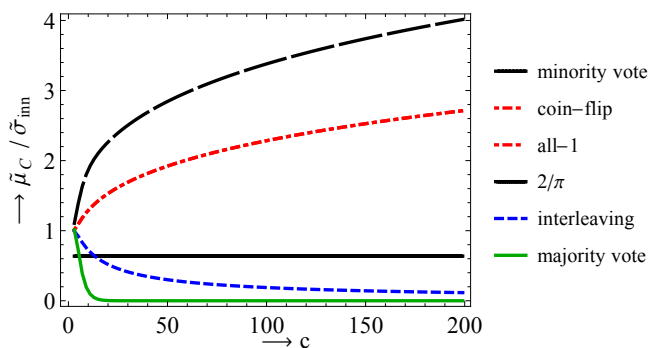


Fig. 4. Coin-flip defense against various attacks in the binary case.

VIII. MAJORITY VOTING DEFENSE

A. Optimal defense

The majority voting attack outputs the symbol with the highest tally. In case of a tie, a uniform choice is made from

the winners. For the binary case, this can be expressed as:

$$f_{y|\mathbf{m}} = \begin{cases} 1 & \text{if } m_y > \frac{1}{2}c \\ \frac{1}{2} & \text{if } m_y = \frac{1}{2}c \\ 0 & \text{if } m_y < \frac{1}{2}c. \end{cases} \quad (170)$$

Lemma 31. Let $q = 2$. For the majority voting attack, we find

$$f_{y|\mathbf{p}}|\mathbf{p}|^c = \sum_{m_y=(c+1)/2}^c \binom{c}{m_y} p_y^{m_y} p_{1-y}^{c-m_y} \quad (171)$$

if c is odd and

$$f_{y|\mathbf{p}}|\mathbf{p}|^c = \frac{1}{2} \binom{c}{c/2} (p_y p_{1-y})^{c/2} + \sum_{m_y=(c+2)/2}^c \binom{c}{m_y} p_y^{m_y} p_{1-y}^{c-m_y} \quad (172)$$

if c is even.

Proof: If c is odd, then

$$f_{y|\mathbf{p}} = \mathbb{E}_{M|\mathbf{p}} [f_{y|M}] = \frac{1}{|\mathbf{p}|^c} \sum_{m_y=\lfloor c/2 \rfloor + 1}^c \binom{c}{m_y} p_y^{m_y} p_{1-y}^{c-m_y}. \quad (173)$$

If instead c is even, the expression receives an additional term $\frac{1}{2} \binom{c}{c/2} (p_y p_{1-y})^{c/2}$. ■

B. Majority Voting Defense Numerics

We ran simulations for the binary case and the arcsine distribution (without cut-off) with the same parameters as described in Section V-D. At first glance, the results look even more promising than those from the interleaving defense. Except against the interleaving attack, the performance of the majority voting defense grows as $c^{0.25}$ against the other 4 considered attacks. However, capacity is not achieved against the interleaving attack.

We note again that the all-1 and coin-flip attacks numerically perform the same against this defense. However, we were unable to show this analytically as we could for the interleaving defense.

We again stress that these five attacks are by no means exhaustive.

TABLE IV

Numerical trends for the performance indicator $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}}$ of the majority voting defense in the binary case $q = 2$ for large c .

interleaving attack	0.91
all-1 attack	$0.66c^{0.22}$
coin-flip attack	$0.66c^{0.22}$
majority voting attack	$0.77c^{0.25}$
minority voting attack	$0.90c^{0.23}$

IX. MINORITY VOTING DEFENSE

A. Optimal defense

The minority voting attack outputs the symbol with the lowest *nonzero* tally. In case of a tie, a uniform choice is made

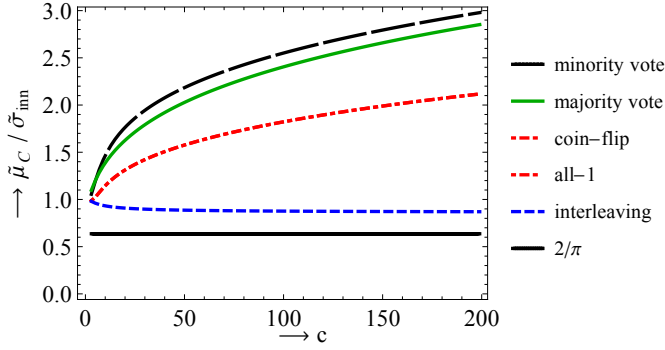


Fig. 5. Majority voting defense against various attacks in the binary case.

from the winners. For the binary case, this can be expressed as:

$$f_{y|m} = \begin{cases} 1 & \text{if } 0 < m_y < \frac{1}{2}c \text{ or } m_y = c \\ \frac{1}{2} & \text{if } m_y = \frac{1}{2}c \\ 0 & \text{if } m_y = 0 \text{ or } \frac{1}{2}c < m_y < c. \end{cases} \quad (174)$$

Lemma 32. Let $q = 2$. For the minority voting attack, we find

$$f_{y|\mathbf{p}}|\mathbf{p}|^c = p_y^c + \sum_{m_y=1}^{(c-1)/2} \binom{c}{m_y} p_y^{m_y} p_{1-y}^{c-m_y} \quad (175)$$

if c is odd and

$$f_{y|\mathbf{p}}|\mathbf{p}|^c = \frac{1}{2} \binom{c}{c/2} (p_y p_{1-y})^{c/2} + p_y^c + \sum_{m_y=1}^{(c-2)/2} \binom{c}{m_y} p_y^{m_y} p_{1-y}^{c-m_y} \quad (176)$$

if c is even.

Proof: If c is odd, then

$$f_{y|\mathbf{p}} = \mathbb{E}_{M|\mathbf{p}} [f_{y|M}] = \frac{1}{|\mathbf{p}|^c} \left(p_y^c + \sum_{m_y=1}^{\lceil c/2 \rceil - 1} \binom{c}{m_y} p_y^{m_y} p_{1-y}^{c-m_y} \right). \quad (177)$$

If instead c is even, the expression receives an additional term $\frac{1}{2} \binom{c}{c/2} (p_y p_{1-y})^{c/2}$.

B. Minority Voting Defense Numerics

We ran simulations for the binary case and the arcsine distribution (without cut-off) with the same parameters as described in Section V-D. The performance against the (targeted) minority voting attack is excellent, and in fact the best when one considers each attack against the matching defense. However, against the other four considered attacks the performance is poor: the minority voting defense fails against the interleaving and majority voting attacks, and only attains a small advantage against the all-1 and coin-flip attacks.

We again see that the all-1 and coin-flip attacks numerically perform the same against this defense. However, we were unable to prove this as we could for the interleaving defense.

We again stress that these five attacks are by no means exhaustive.

TABLE V
Numerical trends for the performance indicator $\tilde{\mu}_c / \tilde{\sigma}_{\text{inn}}$ of the minority voting defense in the binary case $q = 2$ for large c .

interleaving attack	-0.08
all-1 attack	$3.2c^{-0.51}$
coin-flip attack	$3.2c^{-0.51}$
majority voting attack	$-1.9c^{-0.52}$
minority voting attack	$1.4c^{0.25}$

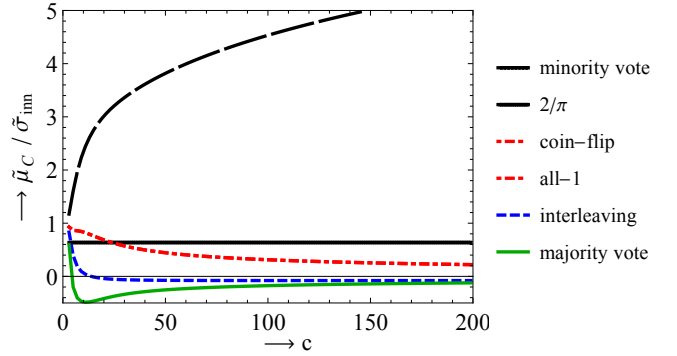


Fig. 6. Minority voting defense against various attacks in the binary case.

X. TARDOS SUSPICION FUNCTION

We end by analyzing the performance of the traditional symmetric Tardos suspicion function.

Lemma 33. If the tracer uses the symmetric Tardos suspicion function, then

$$\tilde{\mu}_c = c \mathbb{E}_{\mathbf{P}} \mathbb{E}_{Y|\mathbf{P}} \left[P_Y \left(\sqrt{\frac{1-P_Y}{P_Y}} - \sqrt{\frac{P_Y}{1-P_Y}} \right) T(Y, Y, \mathbf{P}) - \sqrt{\frac{P_Y}{1-P_Y}} \right]. \quad (178)$$

■ *Proof:* See (3). Since, for fixed y , $h(x, y, \mathbf{p})$ is the same for all $x \neq y$, we find

$$\tilde{\mu}_c = c \cdot \mathbb{E}[T \cdot h] \quad (179)$$

$$= c \mathbb{E}_{\mathbf{P}} \mathbb{E}_{Y|\mathbf{P}} \left[P_Y \sqrt{\frac{1-P_Y}{P_Y}} T(Y, Y, \mathbf{P}) \right. \quad (180)$$

$$\left. - \sqrt{\frac{P_Y}{1-P_Y}} \sum_{x \neq Y} P_x T(X, Y, \mathbf{P}) \right]$$

$$= c \mathbb{E}_{\mathbf{P}} \mathbb{E}_{Y|\mathbf{P}} \left[P_Y \left(\sqrt{\frac{1-P_Y}{P_Y}} - \sqrt{\frac{P_Y}{1-P_Y}} \right) T(Y, Y, \mathbf{P}) \right. \quad (181)$$

$$\left. - \sqrt{\frac{P_Y}{1-P_Y}} \right]$$

■

Against the interleaving attack, the symmetric Tardos suspicion function does not perform well for large q :

Proposition 34. *If the tracer uses the symmetric Tardos suspicion function and the coalition uses the interleaving attack, then $\tilde{\mu}_C = \sum_{y \in \mathcal{A}} \mathbb{E}_{\mathbf{P}}[\sqrt{P_y(1-P_y)}]$. When $f_{\mathbf{P}}$ is a symmetric Dirichlet distribution with concentration parameter $\kappa = \frac{1}{q}$ and no cutoff is used,*

$$\tilde{\mu}_C = \begin{cases} \frac{2}{\pi} & \text{for } q = 2 \\ \frac{1}{2}(q-2) \tan\left(\frac{\pi}{q}\right) & \text{for } q > 2 \\ \frac{\pi}{2} & \text{as } q \rightarrow \infty \end{cases} \quad (182)$$

Proof: When $q = 2$ and p_1 follows the arcsine distribution on $[\delta, 1 - \delta]$ with probability density function (2) then

$$\tilde{\mu}_C = 2 \cdot \mathbb{E}_{\mathbf{P}} \sqrt{p_1(1-p_1)} = \frac{1-2\delta}{\arcsin(1-2\delta)}. \quad (183)$$

For $\delta = 0$ we find $\tilde{\mu}_C = \frac{2}{\pi}$.

Since the marginal distribution of the symmetric Dirichlet distribution is the Beta distribution with parameters κ and $(q-1)\kappa$, we find:

$$\tilde{\mu}_C = \sum_{y=1}^q \mathbb{E}_{\mathbf{P}} \sqrt{p_y(1-p_y)} \quad (184)$$

$$= \sum_{y=1}^q \frac{1}{B(\kappa, (q-1)\kappa)} \int_0^1 p_y^{\kappa+\frac{1}{2}-1} (1-p_y)^{(q-1)\kappa+\frac{1}{2}-1} dp_y \quad (185)$$

$$= q \frac{B(\kappa + \frac{1}{2}, (q-1)\kappa + \frac{1}{2})}{B(\kappa, (q-1)\kappa)} \quad (186)$$

$$= q \frac{\Gamma(\kappa + \frac{1}{2})\Gamma[(q-1)\kappa + \frac{1}{2}]\Gamma(\kappa q)}{\Gamma(q\kappa + 1)\Gamma(\kappa)\Gamma[(q-1)\kappa]} \quad (187)$$

$$= \frac{1}{\kappa} \cdot \frac{\Gamma(\kappa + \frac{1}{2})\Gamma[(q-1)\kappa + \frac{1}{2}]}{\Gamma(\kappa)\Gamma[(q-1)\kappa]}. \quad (188)$$

Now we set $\kappa = \frac{1}{q}$. Using Euler's reflection formula $\Gamma(z)\Gamma(1-z) = \frac{\pi}{\sin(\pi z)}$, we find

$$\tilde{\mu}_C = q \cdot \frac{\sin(\frac{\pi}{q})}{\pi} \cdot \Gamma\left(\frac{1}{q} + \frac{1}{2}\right)\Gamma\left[1 - \frac{1}{q} + \frac{1}{2}\right] \quad (189)$$

$$= q \cdot \frac{\sin(\frac{\pi}{q})}{\pi} \cdot \left(\frac{1}{q} - \frac{1}{2}\right) \cdot \Gamma\left(\frac{1}{q} - \frac{1}{2}\right)\Gamma\left[1 - \frac{1}{q} + \frac{1}{2}\right] \quad (190)$$

$$= \left(1 - \frac{q}{2}\right) \cdot \frac{\sin(\frac{\pi}{q})}{\sin\left[\left(\frac{1}{q} - \frac{1}{2}\right)\pi\right]} = \frac{1}{2}(q-2) \tan\left(\frac{\pi}{q}\right). \quad (191)$$

We see that $\tilde{\mu}_C$ is only a slowly increasing function of q approaching the constant value $\pi/2$, which is far from the optimal code rate.

Proposition 35. *If the tracer uses the symmetric Tardos suspicion function and the coalition uses the all-high attack,*

then

$$\tilde{\mu}_C = c \sum_{y=0}^{q-1} \mathbb{E}_{\mathbf{P}} \left[P_y \left(\sqrt{\frac{1-P_y}{P_y}} - \sqrt{\frac{P_y}{1-P_y}} \right) A_{y+1}^{c-1} \right. \\ \left. - \sqrt{\frac{P_y}{1-P_y}} (A_{y+1}^c - A_y^c) \right]. \quad (192)$$

Proof: This follows directly from Lemma 33 with (113) and (118). ■

Proposition 36. *If the tracer uses the symmetric Tardos suspicion function and the coalition uses the random-symbol attack, then*

$$\tilde{\mu}_C = c \sum_{y=0}^{q-1} \mathbb{E}_{\mathbf{P}} \left[P_y \left[\sqrt{\frac{1-P_y}{P_y}} - \sqrt{\frac{P_y}{1-P_y}} \right] \left[\frac{1}{q} + \sum_{\mathcal{B} \subset \mathcal{A}: y \in \mathcal{B}} \frac{a_{\mathcal{B}}^{c-1}}{|\mathcal{B}|(|\mathcal{B}|+1)} \right] \right. \\ \left. - \sqrt{\frac{P_y}{1-P_y}} \left(\frac{1 - (1-P_y)^c}{q} + \sum_{\mathcal{B} \subset \mathcal{A}: y \in \mathcal{B}} \frac{a_{\mathcal{B}}^c - a_{\mathcal{B} \setminus \{y\}}^c}{|\mathcal{B}|(|\mathcal{B}|+1)} \right) \right]. \quad (193)$$

Proof: This follows directly from Lemma 33 with (155) and (157). ■

It is already known that in the binary case the Tardos defense has a constant $\tilde{\mu}_C$:

Proposition 37. [15] *Let $q = 2$ and $f_{\mathbf{P}}$ be the symmetric Dirichlet distribution with parameter $\kappa = \frac{1}{2}$ without cutoff. If the tracer uses the symmetric Tardos defense, then $\tilde{\mu}_C = \frac{2}{\pi}$, no matter what attack the coalition uses.*

XI. DISCUSSION

We have investigated the optimization of the performance indicator $\tilde{\mu}_C/\tilde{\sigma}_{\text{inn}}$ for bias-based traitor tracing in the simple-decoder setting. A straightforward Lagrangian approach yields a simple expression (Theorem 1) for the optimal suspicion function in a wide variety of contexts, e.g. CDM and RDM, binary and q -ary. The result is a Neyman-Pearson score for the hypothesis $j \in \mathcal{C}$ based on single-location information. It also has the form of a Fisher score, though without a fully understood interpretation.

The h function we obtain with the Lagrangian method depends either on the collusion strategy or on the coalition's symbol tallies \mathbf{m} . These quantities are usually unknown to the tracer. Our optimization approach does not allow for deriving suspicion functions that are based purely on data known to the tracer.

In Section III-A we speculated on the use of the \mathbf{m} -dependent suspicion function in the EM algorithm or as a consistency check for candidate coalitions. Further exploration is left for future work.

TABLE VI
Numerical trends for the performance indicator $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}}$ in the binary case $q = 2$ for large c .

	interleaving attack	all-1 attack	coin-flip attack	majority voting attack	minority voting attack
Tardos defense	$2/\pi$	$2/\pi$	$2/\pi$	$2/\pi$	$2/\pi$
interleaving defense	1.0	$0.61c^{0.23}$	$0.61c^{0.23}$	1.2	$0.75c^{0.25}$
all-1 defense	0.71	$0.86c^{0.25}$	$0.44c^{0.23}$	0.84	$0.54c^{0.25}$
coin-flip defense	$5.1c^{-0.71}$	$0.72c^{0.25}$	$0.72c^{0.25}$	0.0	$1.1c^{0.25}$
majority voting defense	0.91	$0.66c^{0.22}$	$0.66c^{0.22}$	$0.77c^{0.25}$	$0.90c^{0.23}$
minority voting defense	-0.08	$3.2c^{-0.51}$	$3.2c^{-0.51}$	$-1.9c^{-0.52}$	$1.4c^{0.25}$

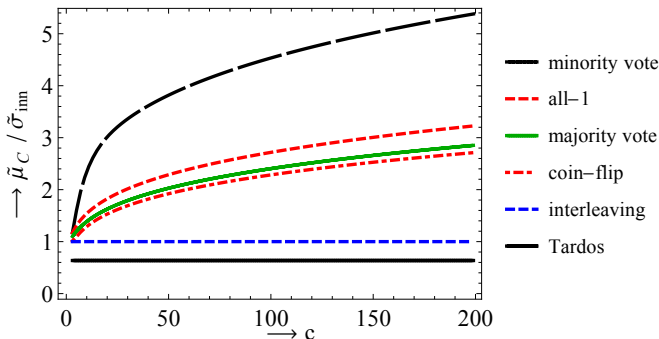


Fig. 7. Performance of optimal suspicion functions against the corresponding attack in the binary case.

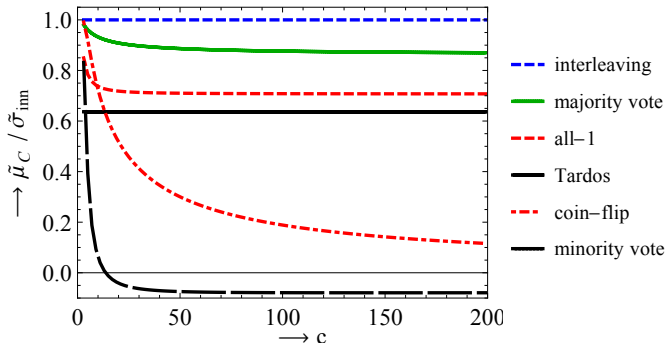


Fig. 8. Interleaving attack against various defenses in the binary case.

For several binary and q -ary attacks in the RDM we have derived the optimal suspicion function. We have investigated the performance indicator $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}}$ in many combinations of suspicion function and attack strategy. In some cases analytic results are obtained. Notably, the matching case of the q -ary interleaving attack gives $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}} = \tilde{\mu}_c = \sqrt{q-1}$, asymptotically ($c \rightarrow \infty$) yielding a code rate precisely equal to the channel capacity [19].

For $q = 2$ the numerical results for the performance indicator $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}}$ are summarized in Table VI. We observe that the interleaving defense, all-1 defense and majority voting defense outperform the Tardos suspicion function for all the considered attacks. In many cases even a positive power of c occurs instead of a constant value: $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}} \propto c^0$ changes to $c^{1/4}$. This is a huge reduction as it leads to a codeword length of $\ell \propto c^{3/2}$. Figure 7 depicts the performance of the optimal defenses against the corresponding attacks. This figure shows that the interleaving attack is particularly strong, as it is the

only one with a constant value of $\tilde{\mu}_c$. The other attacks all seem to scale as $c^{1/4}$, with minority voting being the attack easiest to defend against. Figure 8 shows the performance of the interleaving defense against the five considered attacks.

Another intriguing pattern from the numerical data is the similarity of the all-1 and coin-flip attacks. Except against the all-1 defense, they have the exact same numerical results. Even though for the all-1 attack against the coin-flip defense $\tilde{\sigma}_{\text{inn}} \neq 1$, the normalized $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}}$ values are again the same. We have proven this against the interleaving defense in Proposition 16. This similarity can be explained by realizing that after the collusion attack is performed, the tracer can flip all symbols in the locations where the coalition produced a 0. This transforms the coin-flip attack into the all-1 attack, with the caveat that the coalition then never can receive the $\mathbf{0}$ vector. Naturally, this does NOT apply to the all-1 defense, as this score function is not symbol-symmetric.

It is dangerous to draw general conclusions from the table, however, since not all possible attacks are listed.

Proposition 17 and Theorem 2 on the other hand represent a very important general large- c result: the point (interleaving attack, Dirichlet bias distribution with $\kappa = 1/2$) is a saddle-point of the $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}}$ minimax game when the interleaving defense is used. With the interleaving defense as the simple decoder, the attackers cannot mount a stronger attack than the interleaving attack, and even then they cannot push the rate below the capacity.

With perfect hindsight, this result should not surprise us too much. It was shown by Huang and Moulin [20] that, in the large- c limit, the joint-decoder capacity and simple-decoder capacity coincide. Thus, asymptotically, an optimal simple decoder should automatically achieve capacity.

We now might simply decide to completely switch to the interleaving defense and abandon all other simple decoders. However, the results of Sections IV–X suggest that the other defenses can be used advantageously in a practical decoder scheme at non-asymptotic c . We envisage a decoder that runs the interleaving defense and a small battery of our h functions in parallel (one for every known ‘basic’ strategy, e.g. the ones discussed in this paper). Whenever the colluders use one of the basic strategies, the associated h function will quickly distinguish them from the innocent users; for other strategies, the interleaving defense does the job. The challenge is to combine the different score systems into an effective decoder. Here it has to be borne in mind that both the computational load and the total false positive probability grow with the

number of incorporated h functions.

Future work will focus on (a) investigating which (if any) cutoff to use in a practical traitor tracing scheme, as we expect its scaling behaviour to change; (b) more accurate estimations of σ_{inn} in a practical traitor tracing scheme; (c) efficiency of the interleaving defense at small c , i.e. the non-asymptotic regime where $\tilde{\mu}_c/\tilde{\sigma}_{\text{inn}}$ is no longer the right performance indicator; (d) simulations using multiple suspicion functions in parallel; (e) iterative joint decoders employing the m -dependent suspicion functions.

ACKNOWLEDGEMENTS

We would like to thank Dion Boesten, Thijs Laarhoven, Antonino Simone and Benne de Weger for valuable discussions and comments.

REFERENCES

- [1] G. Tardos, "Optimal Probabilistic Fingerprint Codes," in *Proceedings of the Thirty-Fifth Annual ACM Symposium on Theory of Computing (STOC '03)*, 2003, pp. 116–125.
- [2] O. Blayer and T. Tassa, "Improved versions of Tardos' fingerprinting scheme," *Designs, Codes and Cryptography*, vol. 48, no. 1, pp. 79–103, 2008.
- [3] T. Furon, A. Guyader, and F. C erou, "On the design and optimization of Tardos probabilistic fingerprinting codes," in *Information Hiding*, ser. Lecture Notes in Computer Science, vol. 5284. Springer, 2008, pp. 341–356.
- [4] T. Furon, L. P erez-Freire, A. Guyader, and F. C erou, "Estimating the minimal length of Tardos code," in *Information Hiding*, ser. LNCS, vol. 5806, 2009, pp. 176–190.
- [5] T. Laarhoven and B. de Weger, "Optimal symmetric Tardos traitor tracing schemes," *Designs, Codes and Cryptography*, pp. 1–21, 2012.
- [6] A. Simone and B. Škorić, "Accusation probabilities in Tardos codes: beyond the Gaussian approximation," *Designs, Codes and Cryptography*, vol. 63, no. 3, pp. 379–412, 2012.
- [7] B. Škorić, T. U. Vladimirova, M. U. Celik, and J. C. Talstra, "Tardos Fingerprinting is Better Than We Thought," *IEEE Transactions on Information Theory*, vol. 54, no. 8, pp. 3663–3676, 2008.
- [8] Y.-W. Huang and P. Moulin, "Capacity-achieving fingerprint decoding," in *IEEE Workshop on Information Forensics and Security*, 2009, pp. 51–55.
- [9] K. Nuida, "Short collusion-secure fingerprint codes against three pirates," in *Information Hiding*, ser. LNCS, vol. 6387. Springer, 2010, pp. 86–102.
- [10] K. Nuida, S. Fujitsu, M. Hagiwara, T. Kitagawa, H. Watanabe, K. Ogawa, and H. Imai, "An improvement of discrete Tardos fingerprinting codes," *Des. Codes Cryptography*, vol. 52, no. 3, pp. 339–362, 2009.
- [11] E. Amiri and G. Tardos, "High rate fingerprinting codes and the fingerprinting capacity," in *Proc. 20th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2009, pp. 336–345.
- [12] A. Charpentier, F. Xie, C. Fontaine, and T. Furon, "Expectation maximization decoding of Tardos probabilistic fingerprinting code," in *Media Forensics and Security*, ser. SPIE Proceedings, vol. 7254, 2009, p. 72540.
- [13] P. Meerwald and T. Furon, "Towards joint Tardos decoding: the 'Don Quixote' algorithm," in *Information Hiding*, ser. LNCS, vol. 6958. Springer, 2011, pp. 28–42.
- [14] A. Charpentier, C. Fontaine, T. Furon, and I. Cox, "An asymmetric fingerprinting scheme based on Tardos codes," in *Information Hiding*, ser. LNCS, vol. 6958. Springer, 2011, pp. 43–58.
- [15] B. Škorić, S. Katzenbeisser, and M. U. Celik, "Symmetric Tardos Fingerprinting Codes for Arbitrary Alphabet Sizes," *Designs, Codes and Cryptography*, vol. 46, no. 2, pp. 137–166, 2008.
- [16] B. Škorić, S. Katzenbeisser, H. G. Schaathun, and M. U. Celik, "Tardos Fingerprinting Codes in the Combined Digit Model," *IEEE Transactions on Information Forensics and Security*, vol. 6, no. 3, pp. 906–919, 2011.
- [17] F. Xie, T. Furon, and C. Fontaine, "On-off keying modulation and Tardos fingerprinting," in *Proc. 10th Workshop on Multimedia & Security (MM&Sec)*. ACM, 2008, pp. 101–106.
- [18] B. Škorić and J.-J. Oosterwijk, "Binary and q -ary Tardos codes, revisited," *Cryptology ePrint Archive*, Report 2012/249, 2012.
- [19] D. Boesten and B. Škorić, "Asymptotic fingerprinting capacity for non-binary alphabets," in *Information Hiding 2011*, ser. LNCS, vol. 6958. Springer, 2011, pp. 1–13.
- [20] Y.-W. Huang and P. Moulin, "On fingerprinting capacity games for arbitrary alphabets and their asymptotics," in *IEEE International Symposium on Information Theory (ISIT)*, July 2012, pp. 2571–2575.
- [21] T. Laarhoven, J.-J. Oosterwijk, and J. Doumen, "Dynamic traitor tracing for arbitrary alphabets: Divide and conquer," in *Information Forensics and Security (WIFS), 2012 IEEE International Workshop on*, dec. 2012, pp. 240–245.
- [22] J.-J. Oosterwijk, B. Škorić, and J. Doumen, "Optimal Suspicion Functions for Tardos Traitor Tracing Schemes," in *IH&MMSEC '13*, 2013.