

# Blockwise $p$ -Tampering Attacks on Cryptographic Primitives, Extractors, and Learners

Saeed Mahloujifar\*

Mohammad Mahmoody†

## Abstract

Austrin, Chung, Mahmoody, Pass and Seth [1] studied the notion of bitwise  $p$ -tampering attacks over randomized algorithms in which an efficient ‘virus’ gets to control each bit of the randomness with independent probability  $p$  in an online way. The work of [1] showed how to break certain ‘privacy primitives’ (e.g., encryption, commitments, etc.) through bitwise  $p$ -tampering, by giving a bitwise  $p$ -tampering *biasing* attack for increasing the average  $\mathbb{E}[f(U_n)]$  of any efficient function  $f: \{0, 1\}^n \mapsto [-1, +1]$  by  $\Omega(p \cdot \text{Var}[f(U_n)])$  where  $\text{Var}[f(U_n)]$  is the variance of  $f(U_n)$ .

In this work, we revisit and extend the bitwise tampering model of [1] to *blockwise* setting, where blocks of randomness becomes tamperable with independent probability  $p$ . Our main result is an efficient blockwise  $p$ -tampering attack to bias the average  $\mathbb{E}[f(\bar{X})]$  of any efficient function  $f$  mapping arbitrary  $\bar{X}$  to  $[-1, +1]$  by  $\Omega(p \cdot \text{Var}[f(\bar{X})])$  *regardless* of how  $\bar{X}$  is partitioned into individually tamperable blocks  $\bar{X} = (X_1, \dots, X_n)$ . Relying on previous works of [1, 19, 36], our main biasing attack immediately implies efficient attacks against the privacy primitives as well as seedless multi-source extractors, in a model where the attacker gets to tamper with each block (or source) of the randomness with independent probability  $p$ . Further, we show how to increase the classification error of deterministic learners in the so called ‘targeted poisoning’ attack model under Valiant’s adversarial noise. In this model, an attacker has a ‘target’ test data  $d$  in mind and wishes to increase the error of classifying  $d$  while she gets to tamper with each training example with independent probability  $p$  in an online way.

---

\*University of Virginia, saeed@virginia.edu. Supported by University of Virginia’s SEAS Research Innovation Award.

†University of Virginia, mohammad@virginia.edu. Supported by NSF CAREER award CCF-1350939, and University of Virginia’s SEAS Research Innovation Award.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Our Results	4
1.1.1	Attacks on Randomness of Cryptographic Primitives	4
1.1.2	Efficient Attacks for Biasing Extractors	5
1.1.3	Attacking Learners	6
1.2	Ideas behind Our Blockwise $p$ -Tampering Biasing Attack	7
1.3	Further Related Work and Models	9
<b>2</b>	<b>Preliminaries</b>	<b>9</b>
2.1	Distance Measures	10
2.2	Santha-Vazirani Sources and Their Generalizations	11
<b>3</b>	<b>Blockwise <math>p</math>-Tampering: Definitions and Main Results</b>	<b>12</b>
3.1	Main Results: Blockwise $p$ -Tampering of Bounded Functions	14
<b>4</b>	<b>Applications of <math>p</math>-Tampering Biasing Attacks</b>	<b>16</b>
4.1	Efficient $p$ -Tampering Attacks on Extractors	16
4.2	Targeted Poisoning Attacks on Learners	18
<b>5</b>	<b>Efficient <math>p</math>-Tampering Attacks Biasing Bounded Functions</b>	<b>21</b>
5.1	Biasing Real-Output Functions: Proving Theorem 3.9	23
5.2	Biasing Boolean Functions: Proving Theorem 3.11	32
5.2.1	Part 1: Ideal (Inefficient) Greedy Tampering	32
5.2.2	Part 2: Efficient Greedy Tampering	34
<b>6</b>	<b>Open Questions</b>	<b>38</b>
<b>A</b>	<b>Blockwise <math>p</math>-Tampering Attacks on Primitives: Case of Encryption</b>	<b>42</b>
<b>B</b>	<b>Reducing Blockwise Tampering to Bitwise for Uniform Distributions</b>	<b>44</b>
<b>C</b>	<b>Power and Limitation of Inefficient <math>p</math>-Tampering Attacks</b>	<b>45</b>

# 1 Introduction

In this work, we study *tampering* attacks that efficiently manipulate the *randomness* of randomized algorithms with adversarial goals in mind. Tampering attacks could naturally be studied in the context of cryptographic algorithms that (wish to) access perfectly uniform and untampered randomness for sake of achieving security. However, the scope of such attacks goes beyond the context of cryptography and could be studied more broadly for any class of algorithms that depend on some form of untampered random input and try to achieve specific goals (e.g., learning algorithms using untampered training data to generate a hypothesis). Here, we are interested in understanding the power and limitations of such tampering attacks over the randomness, when the adversary can tamper with, or even control,  $\approx p$  fraction of the randomness.<sup>1</sup>

The most relevant to our study here is the work of Austrin et al. [1] that introduced the notion of *bitwise  $p$ -tampering* attacks on the randomness of cryptographic primitives. In this model, the adversary generates an efficient ‘virus’ who gets into the ‘infected’ device, can read everything, but is limited in what it can alter. As the stream of bits of randomness  $R = (r_1, \dots, r_n)$  is being generated, for every bit  $r_i$ , the  $p$ -tampering virus gets to change  $r_i$  with independent probability  $p$  (i.e., with probability  $(1-p)$  the bit remains unchanged).  $p$ -tampering attacks are online, so the virus does not know the future incoming bits, but it can base its decisions based on the history of the (potentially tampered) bits. The work of [1] proved that bitwise  $p$ -tampering attacks can always increase the average of efficient bounded functions  $f: \{0, 1\}^n \mapsto [-1, +1]$  by  $\Omega(p \cdot \text{Var}[f(U_n)])$  where  $\text{Var}[f(U_n)]$  is the variance of  $f(U_n)$ .

Austrin et al. [1] showed how to break a variety of ‘privacy’ cryptographic primitives (e.g., public-key and private key encryption, zero knowledge, commitments, etc.) that have ‘indistinguishability-based’ security games using their main efficient bitwise  $p$ -tampering biasing attack. In such cryptographic attacks, the *code* of the  $p$ -tampering virus is generated by an outside adversary who only knows the public information (e.g. public key). Previously, Dodis, Ong, Prabhakaran, and Sahai [19] had shown that for the same cryptographic primitives, there are high-min-entropy Santha-Vazirani sources of randomness [39] that make them insecure. Thus, the work of [1] was a strengthening of the results of [19] showing how to generate such ‘bad’ SV sources through efficient  $p$ -tampering attacks. The  $p$ -tampering attacks of [1], and in particular their core attack for biasing the output of balanced bounded functions, crucially depend on the fact that the attacker can tamper with *every single bit* of the randomness *independently* with probability  $p$ . However, randomness is usually generated in blocks rather than bits [4, 16, 21, 28], e.g., during the boot time [30], and is also made available to the algorithms requesting them in blocks. Thus, it is indeed natural to consider tampering attackers who sometimes get to change an incoming *block* of randomness.

**Blockwise  $p$ -tampering attacks.** In this work, we revisit the bitwise  $p$ -tampering model of [1] and extend it to a setting where the tampering could happen over blocks. Suppose  $A$  is an algorithm taking  $\bar{X} = (X_1 \times \dots \times X_n)$  as input where  $\bar{X}$  is a distribution consisting of  $n$  blocks and the  $i$ ’th block is independently sampled from the distribution  $X_i$ . For example,  $A$  could be a cryptographic algorithm in which  $X_i$  is the  $i$ ’th block of *uniform* randomness given to  $A$ . Or  $A$  could also be a learning algorithm given  $n$  i.i.d training examples. Roughly speaking, a blockwise  $p$ -tampering attack on (the randomness of)  $A$  is an algorithm Tam working as follows. Suppose we sample the blocks  $x_i \leftarrow X_i$  one by one. Then the  $i$ ’th block  $x_i$  becomes ‘tamperable’ with independent probability  $p$  for each  $i$ , and it remains intact with probability  $1 - p$ . In case  $x_i$  becomes tamperable, then Tam could substitute  $x_i$  with another value  $x'_i$  in the support set<sup>2</sup> of  $X_i$  in an

<sup>1</sup>Note that if the adversary can control all the randomness, we are effectively back to what we can do in the deterministic setting.

<sup>2</sup>We only allow the tampering algorithm to produce something in the support set. A more general definition allows the tampering algorithm to make choices out of the support set, however, our restriction only makes our attacks stronger.

online way. Namely, when Tam gets the chance to tamper with  $x_i$  it could decide on a new block  $x'_i$  based on the knowledge of previous (tampered) blocks. The tampering algorithm Tam could also depend on (and thus know everything about) the algorithm  $A$  including all of its inputs selected so far, but it cannot write anything except when it is given the chance to tamper with a block of randomness.

Different  $p$ -tampering attackers could pursue different goals. For example, as it was done in the bitwise setting of [1], a  $p$ -tampering attack might aim to ‘signal out’ a secret information (e.g., the plain-text). Another example is when Tam wants to increase the classification error of the hypothesis output by a learner  $A$  where each block  $x_i = (d, t)$  consists of a labeled example sampled from the same distribution.

We also note that, though called primarily a tampering attack,  $p$ -tampering attacks are not blind tampering attackers and naturally rely on the knowledge of the previous random bits before deciding on the tampering of the next bit/block, although such knowledge is only given to the tampering virus, and e.g., not the external adversary generating the code of the virus. That is a reason why the proven power of  $p$ -tampering attacks in this work is not in contradiction with known *positive* results such as [18, 24, 26, 32].

## 1.1 Our Results

Our main result is a generalization of the biasing attack of [1] to the blockwise setting. We first describe this result, and then we will describe some of the applications of this biasing attack.

**Theorem 1.1** (Informally stated). *Let  $\bar{X} = (X_1 \times \dots \times X_n)$  be a product distribution where each of  $X_i$ ’s is efficiently samplable. For any efficient function  $f: \text{Supp}(\bar{X}) \mapsto [-1, +1]$  there is an efficient blockwise  $p$ -tampering attack that increases the average of  $f$  over a sampled input by at least  $\Omega(p) \cdot \text{Var}[f(\bar{X})]$ .*

See Theorem 3.9 for a formalization. Similarly to [1], we also prove a variant of Theorem 1.1 for the special case of Boolean functions, but with better parameters (see Theorem 3.11). However, some of the applications of this biasing lemma (e.g., for attacking cryptographic primitives, or attacking learning algorithms with non-Boolean cost/loss functions) we need to use the non-Boolean attack of Theorem 1.1.

Our main biasing  $p$ -tampering attack on bounded functions even applies to the settings where  $\bar{X}$  is *not* a product distribution. In that case, we assume that  $\bar{X}$  is sampled in a ‘stateful’ way, and that the next block  $X_i$  is sampled conditioned on adversary’s choices of blocks. This extension allows our model to include previous special models of  $p$ -tampering attacks against random walks on graphs [3].

We also prove some applications for our main biasing attack that rely on the *blockwise* nature of it. In addition to obtaining attacks against the security of cryptographic primitives as well as multi-source randomness extractors through blockwise  $p$ -tampering, we also demonstrate applications beyond cryptography. In particular, by relying on the power of biasing attacks over *non-uniform* distributions, we show how to attack and increase the error of *learning* algorithms that output classifiers, through an attack that injects a  $p$  fraction of adversarial data in an online way. In what follows we briefly discuss each of these applications.

### 1.1.1 Attacks on Randomness of Cryptographic Primitives

As mentioned, the bitwise  $p$ -tampering attack of [1] for biasing functions was at the core of their attacks breaking the security of cryptographic primitives by tampering with their randomness. By using our biasing attack of Theorem 1.1 we immediately obtain blockwise attacks against the same primitives. This time, our attacks work *regardless* of how randomness is packed into blocks, and is also ‘robust’ in the sense that the attack succeeds even if the tampering probabilities  $p_1, p_2, \dots$  are *not* equal so long as  $p \leq p_i$  for all  $i$ .<sup>3</sup>

---

<sup>3</sup>In fact, we observe that the bitwise  $p$ -tampering attack of [1] can also be shown to be robust. Moreover, we believe robustness is an important feature for cryptographic attacks and so worth to be studied explicitly.

**Corollary 1.2 (Informal).** *Let  $\mathcal{P}$  be one of the following primitives. CPA secure public-key or private-key encryption, efficient-prover zero-knowledge proofs for  $\mathbf{NP}$ , commitment schemes, or two party computation where only one party gets the output. Then there is an efficient blockwise  $p$ -tampering attack that breaks the security of  $\mathcal{P}$  with advantage  $\Omega(p)$ . In particular, the attack succeeds even if the length of the tampered randomness blocks are unknown a priori and only become clear during the attack.*

The above theorem could be obtained by plugging in our biasing attack of Theorem 1.1 into the proofs of [1]. However, for sake of completeness, in Section A we give a formal definition of robust blockwise attacks breaking the CPA security of public-key encryption schemes, and here we focus on further new applications of our blockwise  $p$ -tampering biasing attacks.

**Achieving security against blockwise  $p$ -tampering?** In addition to presenting the power of bitwise  $p$ -tampering attacks, the work of [1] also showed how to achieve secure protocols against bitwise  $p$ -tampering attacks for ‘forging-based’ primitives such as signatures for  $p = 1/\text{poly}(\kappa)$  where  $\kappa$  is the security parameter. For the same primitives, when we move to the blockwise setting, whether or not achieving positive (secure) results is possible *depends* on the block sizes of the tampering attack. For example, if the *whole* randomness of the key generation algorithm of a signature scheme becomes tamperable as a single block (with probability  $p \geq 1/\text{poly}(\kappa)$ ) the adversary can choose an insecure key. On the other hand, if all the blocks are of constant size (or even of size  $o(\lg \kappa)$ ) similar arguments to those in [1] could be used to make ‘forging-based’ primitives secure for any  $p \leq \kappa^{-\Omega(1)}$ .

### 1.1.2 Efficient Attacks for Biasing Extractors

Our blockwise  $p$ -tampering attacks for biasing functions are natural tools for ‘biasing attacks’ against (seedless) randomness extractors from block sources.

**Biasing multi-source seedless extractors.** We can directly use our  $p$ -tampering attacks against *any* specific, multi-source, seedless randomness extractors [12, 39, 44]. Namely, suppose  $f$  is an efficient seedless extractor who takes  $n$  blocks of randomness  $(x_1, \dots, x_n) \leftarrow (X_1 \times \dots \times X_n)$  where the distribution  $X_i$  belongs to a class of randomness source. Then, for any choice of samplable  $\bar{X} = (X_1, \dots, X_n)$ , Theorem 3.11 gives an efficient  $p$ -tampering attacker who could transform the distribution  $\bar{X}$  into  $\bar{Y}$  such that  $|\mathbb{E}[f(\bar{Y})]| \geq \Omega(p)$ . Note that the interesting aspect of  $\bar{Y}$  is that it is identical to  $\bar{X}$  in  $(\approx 1 - p)$  fraction of the blocks. In particular, as we will see, our attacker of Theorem 1.1 has the property that upon tampering with each block, all it does is to either leave as is or ‘resample’ it once.

The second application of our  $p$ -tampering attacks against extractors is different in the sense that instead of attacking extractors when unbiased extraction is possible, it gives an alternative algorithmic proof for a known impossibility result [6, 19, 22, 36] regarding block Santha-Vazirani sources [39]. Below, by  $U_i^j = U_i \times \dots \times U_i$  we refer to  $j$  blocks each consisting of  $i$  uniform bits.

**Impossibility of randomness extraction from block-SV sources.** The celebrated work of Santha and Vazirani [39] proved a strong negative result about deterministic randomness extraction from sources with high min-entropy. An SV source (see Definition 2.10) is a joint distribution  $(X_1, \dots, X_n)$  over  $\{0, 1\}^n$  with the guarantee that every bit is  $\delta$ -close to uniform even if we condition on all the previous bits. In particular, [39] proved that for any deterministic (supposedly extractor) function  $f: \{0, 1\}^n \mapsto \{+1, -1\}$ , there is always an  $\delta$ -SV source  $\bar{X} = (X_1, \dots, X_n)$  such that  $|\mathbb{E}[f(\bar{X})]| \geq \Omega(\delta)$ . The work of Reingold, Vadhan and Wigderson [36] gave an elegant simple proof for this result using the so called ‘half-space’

sources, and this idea found its way into the work of Dodis et al. [19] where they generalized the result of [39] to *block* sources [13]. A  $(\ell, k)$ -block SV source is a sequence of blocks of length  $\ell$  bits such that each block has min-entropy at least  $k$  conditioned on previous blocks (see Definition 2.11).

Even though  $p$ -tampering attacks do *not* generate block-SV sources with ‘high’ min-entropy in general, we show that the *specific*  $p$ -tampering attacker of our Theorem 1.1 does indeed generate an  $(\ell, \ell - p)$  block-SV source. As a result, we get an alternative proof for the impossibility of deterministic extraction from block-SV sources, but this time through *efficient*  $p$ -tampering attacks.<sup>4</sup> In particular, we prove the following.

**Theorem 1.3** (Efficient  $p$ -tampering attacks against block-SV extractors). *Let the function  $f: \{0, 1\}^{\ell-n} \mapsto \{+1, -1\}$  be a ‘candidate’ efficient deterministic extractor for  $(\ell, \ell - p)$  block SV sources. Then there is an efficient  $p$ -tampering attack that generates a  $(\ell, \ell - p)$  block SV source for which the  $f$  has average  $\Omega(p)$ .*

Our main contribution in Theorem 1.3 is the efficiency of its  $p$ -tampering attacker, as without that condition one can prove Theorem 1.3 using a *computationally unbounded*  $p$ -tampering attacker and the proof implicit in [19, 36] and explicit in [6, 22] for the case of block SV sources. In fact, we prove a more general result than Theorem 1.3 by proving the impossibility of efficient bit bit extractors from yet another generalization of SV sources, called mutual max-divergence [23] (MMD) sources (see Definition 2.8).

### 1.1.3 Attacking Learners

In this work, we also use our blockwise  $p$ -tampering attack in the context of “adversarial” machine learning [5, 35] where an adversary aims to increase the error of a learning algorithms for a specific test data that is known to him. In what follows, the reader might find the review of the standard terminology at the beginning of Section 4.2 useful.

**Targeted poisoning attacks against learners.** Poisoning attacks (a.k.a causative attacks) [2, 41, 45] model attacks against learning systems in which the adversary manipulates the training data  $\bar{x} = (x_1, \dots, x_n)$ , where  $x_i$  is the  $i$ ’th *labeled* training example, in order to increase the error of the learning algorithm. Poisoning attacks could model scenarios where the tampering happens *over time* [37, 38] e.g., because the learning algorithm is retrained daily or weekly using potentially tamperable data. *Targeted* (poisoning) attacks [41] refer to the setting where the adversary *knows* a specific test data  $\mathcal{X}$  over which the hypothesis will be tested, and she probably has some interest in increasing the error of the hypothesis over that particular test set  $\mathcal{X}$ . For simplicity of discussion, below we assume that  $\mathcal{X} = \{(d, t)\}$  where  $t$  is the label of  $d$  and the adversary’s goal is to make the learning algorithm output a wrong label for  $d$ .

A very natural model for how the poisoning attacks occur was defined by Valiant [43]. In this model, a training oracle  $O_X(\cdot)$  for a distribution  $X$  (from which the training sequence  $\bar{x} = (x_1, \dots, x_n)$  will be sampled) would be manipulated by an adversary as follows. Whenever the training algorithm queries this oracle, with probability  $1 - p$  the answer is generated from the original oracle  $O_X$  and with probability  $p$  a stateful adversary  $A$  gets control over the oracle and answers with an arbitrary pair  $(d, t)$ . Many subsequent work (e.g., [10, 31]) studied how to make learners secure against such noise but *not* in the targeted setting.

**Valiant’s model vs.  $p$ -tampering.** Valiant’s adversarial model for the training oracle is indeed very similar to our blockwise  $p$ -tampering model except for the fact that in the Valiant’s model, the adversary is allowed to use wrong labels (i.e.,  $x_i = (d, t)$  where  $t$  is *not* the correct label for  $d$ ). However, as we discussed above, our  $p$ -tampering attackers are not allowed to go out of the ‘support set’ of the distribution (see Definition 4.5).

<sup>4</sup>Note that this is indeed a stronger condition than just getting a samplable source. See Remark 3.6.

In this work, we prove the following attack against deterministic learners of classifiers (see Theorem 4.7 for a formalization). One subtle difference between the models is that in Valiant’s model, the adversary knows everything about the *current* state of the learner, while in our model, the adversary knows the history of the blocks. For all of our attacks, all adversary needs is to ‘continue’ the computation done by the learner, and knowing the current state (as in Valiant’s model) allows us to do so, even if the previous blocks are unknown. Therefore, all of our  $p$ -tampering attacks indeed apply in Valiant’s model.

**Theorem 1.4** (Informal–Targeted poisoning attacks against classifiers). *Let  $L$  be a deterministic learning algorithm  $L$  that takes a sequence  $\bar{x} = (x_1, \dots, x_n)$  of i.i.d samples from the same distribution  $X$ , where  $x_i = (d_i, \ell_i)$  and  $\ell_i$  is the label of  $d_i$ . Suppose, without tampering, the probability of  $L$  making a mistake on test example  $d$  is  $\delta$  over the choice of  $x_1, \dots, x_n \leftarrow X$ . Then there exists a  $p$ -tampering attack over the training sequence  $(x_1, \dots, x_n)$  that increases the error for classifying  $d$  to  $\delta' \geq \delta + \Omega_\delta(p)$ . Moreover, if  $X$  is efficiently samplable, the attack is efficient as well.*

Note that the above attacker is a  $p$ -tampering one, meaning it never goes out of the support set of the distribution. In other words, our attacker does *not* use any wrong labels in its adversarial samples! Therefore, our attacks are ‘defensible’ in the sense that what they produce is always a possible legitimate outcome of the honest sampling, so it could not be *proved* in court that the data is not generated honestly! Previous work on poisoning attacks (e.g., see [2, 41, 45]) has studied attacks against *specific* learners, while our result can be applied to *any* learner.

**Comparison with the distribution-independent setting of [10,31].** Previous works of Kearns and Li [31] and Bshouty, Eiron, and Kushilevitz [10] have already proved impossibility of PAC learning in Valiant’s model of adversarial noise. In addition to using wrong label in their attacks (which is not allowed in the  $p$ -tampering model) there is also another distinction between their model and our  $p$ -tampering poisoning attacks. The attacks of the works [10, 31] are in the *distribution-independent* setting, and their negative results heavily rely on the *existence* of some initial distribution that is not PAC learnable under adversarial noise. Our attacks, on the other hand, apply even to the *distribution-specific* setting, where the adversary has no control over the initial distribution, and it can always turn that distribution against the learner.

## 1.2 Ideas behind Our Blockwise $p$ -Tampering Biasing Attack

In this subsection we describe some of the ideas behind the proof of our Theorem 1.1.

**Reduction to bitwise tampering?** Our first observation is that blockwise  $\tilde{p}$ -tampering over *uniformly* distribute blocks  $U_{s_1} \times \dots \times U_{s_n}$  could be reduced to  $p$ -tampering over  $N = \sum_i s_i$  many uniform bits, as long as  $1 - \tilde{p} \leq (1 - p)^{s_i}$  for every  $s_i$ . The idea is that if  $1 - \tilde{p} \leq (1 - p)^{s_i}$ , then the probability of the whole block  $U_{s_i}$  getting tampered with in the blockwise model is at least the probability that *at least* one of the bits are tampered with in the bitwise model. Therefore, a blockwise attacker can ‘emulate’ the bitwise attacker internally. In Section B we formally describe this rather simple reduction.

However, this reduction is imperfect in three aspects. (1) Firstly, to use this reduction we will need to use  $p \approx \tilde{p}/s$  where  $s$  is the maximum length of any block. Therefore, we cannot gain any bias more than  $1/s$  which, in particular, would be at most  $o(1)$  for non-constant block sizes  $s = \omega(1)$ . This prevents us from getting applications (e.g., attacks against extractors) that require large  $\Omega(1)$  bias. (2) Secondly, this reduction only works for blocks that are originally distributed as uniform bits (i.e.,  $U_s$ ), and so it cannot be applied to general non-uniform distributions, which is indeed the setting of our  $p$ -tampering attacks against

learners. (3) Finally, this reduction does not preserve robustness as the  $\tilde{p}$ -tampering algorithm would need to know the *exact* probabilities under which the tampering happens, while in our applications of blockwise tampering to cryptographic primitives robustness we aim for robust attacks that do not depend on this exact knowledge. Because of all this, in this work we aim for a direct attack analyzed in the blockwise regime.

The work of [1] used a so called ‘mild-greedy’ attack for biasing real-valued bounded function in a bitwise  $p$ -tampering attack. Roughly speaking, this attack works as follows. When the tampering happens, the tampering algorithm first picks a random bits  $b'_i$ . Then, using a random continuation  $b'_{i+1}, \dots, b'_n$  it interprets  $s = f(b_1, \dots, b_{i-1}, b'_i, \dots)$  as how good the choice of  $b'_i$  is. Then, using a biased coin based on  $s$ , the tampering algorithm either keeps  $b'_i$  or it flips it to  $1 - b'_i$ . This attack, unfortunately, is tailored of the bitwise setting, as flipping a block is not natural (or even well defined).

**Our new one rejection sampling attack.** In this work propose a new attack for the blockwise setting that is inspired by the mild-greedy attack of [1]. Our attack is not exactly a ‘generalization’ of the mild-greedy attack to the blockwise setting, as even for the case of uniform blocks of one bit, it still differs from the mild-greedy attack, but it is nonetheless inspired by the one-greedy attack and its analysis also uses ideas from the analysis of mild-greedy attack [1]. We call our tampering attack *one rejection sampling*, denoted by ORSam, and it works as follows: given previously chosen blocks  $(y_1, \dots, y_{i-1})$  for  $\bar{X}$  (some of which might be the tampered blocks) the tampering algorithm ORSam first samples  $(y'_i \leftarrow X_i, \dots, y'_n \leftarrow X_n)$  ‘in its head’, then gets  $s = f(y_1, \dots, y_{i-1}, y'_i, \dots, y'_n)$ , and outputs:

$$\begin{cases} \text{Case 1: with probability } \frac{1+s}{2} : & \text{keep } y'_i \\ \text{Case 2: with probability } \frac{1-s}{2} : & \text{use a fresh sample } y''_i \leftarrow X_i. \end{cases}$$

**Why does one-rejection sampling work?** The main challenge is to show that the simple one-rejection sampling attack described above actually achieves bias proportional to the variance. In order to relate the bias to the variance of the function, we first need to define two notations. For every prefix  $x_{\leq i} = x_1, \dots, x_i$  let  $\hat{f}[x_{\leq i}] = \mathbb{E}[f(\bar{X}) | X_1 = x_1, \dots, X_i = x_i]$  to be the average of function  $f$  w.r.t to distribution  $\bar{X}$  conditioned on that prefix. Also let  $g[x_{\leq i}] = \hat{f}[x_{\leq i}] - \hat{f}[x_{\leq i-1}]$  be the change in average of  $f$  (i.e.,  $\hat{f}$ ) when we go from  $x_{\leq i-1}$  to  $x_{\leq i}$ . A straightforward calculation shows that

$$\text{Var}[f(\bar{X})] = \mathbb{E}_{(x_1, \dots, x_n) \leftarrow \bar{X}} \left[ \sum_{i \in [n]} g[x_{\leq i}]^2 \right] = \sum_{i \in [n]} \mathbb{E}_{x_{\leq i} \leftarrow (X_1, \dots, X_i)} [g[x_{\leq i}]^2]. \quad (1)$$

That is simply because the sequence  $(\hat{f}[x_{\leq 0}], \dots, \hat{f}[x_{\leq n}])$  forms a martingale. Suppose  $\bar{Y} = (Y_1, \dots, Y_n)$  is the new distribution after the  $p$ -tampering happens over  $\bar{X}$ . Equation (1) suggests the following natural idea for lower bounding the amount of ‘‘global gain’’ that is achieved for increasing the average  $d = \mathbb{E}[f(\bar{Y})] - \mathbb{E}[f(\bar{X})]$  under the attack’s generated distribution by relating it to the variance  $\text{Var}[f(\bar{X})]$ . In particular, it would suffice to lower bound the ‘‘local gains’’ for average of  $f$  when we apply our *one rejection sampling* with probability  $p$  for a particular block  $i$ , by relating it the term  $\mathbb{E}_{(x_1, \dots, x_n) \leftarrow \bar{X}} [g[x_{\leq i}]^2]$  (for the same fixed  $i$ ). Direct calculation shows that the ‘local gain’ obtained by our one-rejection sampling attack for any prefix  $x_{\leq i}$  is *exactly*  $\frac{p}{2} \cdot \mathbb{E}_{x_{i+1} \leftarrow X_{i+1}} [g[x_{\leq i}, x_{i+1}]^2]$ .

Unfortunately, a subtle point prevents us from using the above argument, because as soon tampering happens, we *deviate* from the original distribution  $\bar{X}$ , and the ‘prefixes’ of the blocks come from a new distribution  $\bar{Y}$  rather than  $\bar{X}$ , so we cannot directly use to Equation (1) to lower bound the local gains by relating them to  $\text{Var}[f(\bar{X})]$ . Nonetheless, it can be shown that a variant of Equation (1) still holds in which, roughly speaking,  $\text{Var}[f(\bar{Y})]$  substitutes  $\text{Var}[f(\bar{X})]$ . Therefore, it would be sufficient to lower



bound  $\text{Var}[f(\bar{Y})]$  based on  $\text{Var}[f(\bar{X})]$ . For this goal, we employ similar ideas to those of [1] to show by induction over  $i$  that at any moment during the attack *either* the average *or* the variance of  $\hat{f}[x_{\leq i}]$  under the *new* tampered distribution  $\bar{Y}$  is large enough. See Section 5 for more details.

### 1.3 Further Related Work and Models

Since the work of Boneh, DeMillo and Lipton [9] it is known that even *random* tampering with computation of certain protocols could lead to devastating attacks. The work of Gennaro, Lysyanskaya, Malkin, Micali, and Rabin [26] initiated a formal study of algorithmic tamper resilience. Along this direction, non-malleable codes, introduced by Dziembowski, Pietrzak, and Wichs [25], become a central tool for preventing tampering attacks on the internal state of an algorithm. More recently, Chandran et al. [11] studied non-malleable codes in the *blockwise* tampering model that bears similarities to our model in this work, though our goals are completely different. Finally, Bellare, Paterson, and Rogaway [7] initiated the study of *algorithm substitution* attacks where a powerful attacker can adversarially substitute components of the algorithm.

**Coin-tossing.** At a high level, our blockwise tampering attacks, specially for biasing Boolean functions, have some conceptual similarities to attacks against coin-tossing protocols [8, 15, 17, 29, 34]. Indeed, both types of attacks aim at biasing a final bit by ‘substituting’ some ‘blocks’. In our setting, the block is the next sampled chunk of randomness, and for coin tossing blocks are maliciously chosen messages to the other party! However, the pattern of tampering in such attacks is one out of two complementing sets (referring to each party’s turns), while in our setting each block becomes tamperable with an independent probability  $p$ .

**Tampering with ‘bounded budget’.** The works of [15, 27, 33] studied the power of related tampering attacks in the blockwise setting where the goal of the adversary is indeed to bias the output of a function. However, in these papers, while the adversary has a ‘limited budget’ of how many times to tamper, it *can choose* when to tamper with a block, while, in our model the adversary will have no control on about  $1 - p$  fraction of the blocks, and he does not get to choose which blocks will be so. The work of Dodis [20] studies a ‘mixture’ of both models where the adversary has a bounded budget that she can use upon choice, but she also gets to tamper ‘randomly’ otherwise.

## 2 Preliminaries

Logarithms are denoted by  $\lg(\cdot)$  and, unless specified otherwise, they are in base 2. By  $a, b \in \mathcal{D}$  we mean that  $a \in \mathcal{D}$  and  $b \in \mathcal{D}$ . For a string  $x \in \{0, 1\}^*$ , by  $|x| = n$  we denote that  $x \in \{0, 1\}^n$ . For a randomized algorithm  $S$ , we only explicitly represent its input and do not represent its randomness and by  $y \leftarrow S(x)$  we denote the process of running  $S(x)$  using fresh randomness and getting  $y$  as output.

**Notation on random variables.** Unless specified otherwise, all of the random variables and distributions in this work are discrete and finite. We use uppercase letters to denote random variables and distributions (e.g.,  $X$ ). For real valued random variable  $X$ , by  $\mathbb{E}[X]$  and  $\text{Var}[X]$ , we mean (in order) the expected value and variance of  $X$ . We usually use the same letter to refer to distributions and random variables sampled from them. By  $\text{Supp}(X) = \{x \mid \Pr[X = x] > 0\}$  we denote the support set of  $X$ . The process of sampling  $x$  from  $X$  is denoted by  $x \leftarrow X$  and  $X \equiv Y$  is used to show that  $X$  and  $Y$  are distributed identically.

By  $U_m$  we denote the random variable uniformly distributed over  $\{0, 1\}^m$ . By  $(X, Y)$  we denote random variables  $X, Y$  that are distributed jointly. By  $(X \times Y)$  we mean  $(X, Y)$  where  $X$  and  $Y$  are independently

sampled from their marginal distribution. For joint random variables  $(X, Y)$  and for any  $y \leftarrow Y$ , by  $(X \mid y)$  we denote the distribution of  $X$  conditioned on  $Y = y$ . By using a random variable like  $X$  in an expected value (or probability) we mean that the expected value (or the probability) is also over  $X$  (e.g.,  $\mathbb{E}[f(X)] = \mathbb{E}_{x \leftarrow X}[f(x)]$  and  $\Pr[f(X) = 1] = \Pr_{x \leftarrow X}[f(x) = 1]$ ). We also use the tradition that the multiple appearances of the same random variable  $X$  in the same phrase refer to identical samples (e.g., it always holds that  $\Pr[X = X] = 1$ ). For a random variable  $D$ , we also use  $D(x)$  to denote  $\Pr[D = x]$ .

**Definition 2.1** (Bit extraction). Let  $\mathcal{X}$  be a set of distributions over a domain  $\mathcal{D}$ . We call a function  $f: \mathcal{D} \mapsto \{+1, -1\}$  an  $\varepsilon$ -extractor for  $\mathcal{X}$  (sources) if for every  $X \in \mathcal{X}$  it holds that  $|\mathbb{E}[f(X)]| \leq \varepsilon$ .

**Definition 2.2.**  $H_\infty(X) = \min_{x \in \text{Supp}(X)} \lg(1/p(x))$  is the *min-entropy* of  $X$ .

**Definition 2.3** (Span of distributions). Let  $\mathcal{X} = \{X_1, \dots, X_k\}$  be a set of distributions over the same domain. For  $\alpha_1 + \dots + \alpha_k = 1$ , by  $X = \sum_{i \in [k]} \alpha_i X_i$  we refer to the distribution  $X$  such that  $\Pr[X = a] = X(a) = \sum_i \alpha_i X_i(a)$ . Namely,  $X$  can be sampled by the following process: first sample  $i \in [k]$  with probability  $\alpha_i$ , then sample  $x \leftarrow X_i$  and output  $x$ . The *span* of distributions in  $\mathcal{X}$  is defined to be the set of all convex combinations of distributions in  $\mathcal{X}$ :  $\text{Span}(\mathcal{X}) = \{X = \sum_{i \in [k]} \alpha_i X_i \mid \sum_{i \in [k]} \alpha_i = 1\}$ .

**Lemma 2.4** (Hoeffding's inequality). Suppose  $A_1, \dots, A_n$  are i.i.d random variables distributed over  $[-1, +1]$  with expected value  $\mathbb{E}[A_i] = \mu$ , and let  $A = \mathbb{E}_{i \leftarrow [n]}[A_i]$  be their average. Then, for all  $\varepsilon \geq 0$  we have  $\Pr[|A - \mu| \geq \varepsilon] \leq e^{-n \cdot \varepsilon^2 / 2}$ .

## 2.1 Distance Measures

**Definition 2.5** (Statistical distance). The *statistical distance* (a.k.a. total variation distance) between random variables  $X, Y$  is defined as

$$D_{\text{SD}}(X, Y) = \max_{E \subseteq \text{Supp}(X)} \Pr[X \in E] - \Pr[Y \in E].$$

The following lemma gives a well known characterization of the statistical distance.

**Lemma 2.6** (Characterizing statistical distance). It holds that  $D_{\text{SD}}(X, Y) \leq p$  iff there are distributions  $Z, X', Y'$  such that  $X = (1 - p)Z + pX'$  and  $Y = (1 - p)Z + pY'$ . In particular, if  $Y = (1 - p)X + pZ$  then we have  $D_{\text{SD}}(X, Y) \leq p$  because it always holds that  $X = (1 - p)X + pX$ .

**Definition 2.7** (KL-divergence). The *Kullback-Leibler (KL) divergence* from distribution  $Q$  to distribution  $P$  is defined as follows:  $D_{\text{KL}}(P||Q) = \mathbb{E}_{a \leftarrow P} \lg(P(a)/Q(a))$  if  $\text{Supp}(P) \subseteq \text{Supp}(Q)$ , and  $D_{\text{KL}}(P||Q) = \infty$  if  $\text{Supp}(P) \not\subseteq \text{Supp}(Q)$ .

**Definition 2.8** (Max-divergence [23]). The *max-divergence*  $D_\infty(P||Q)$  from random variable  $Q$  to  $P$  is  $\max_{a \in \text{Supp}(P)} \lg(P(a)/Q(a))$  if  $\text{Supp}(P) \subseteq \text{Supp}(Q)$ , and  $D_\infty(P||Q) = \infty$ , if  $\text{Supp}(P) \not\subseteq \text{Supp}(Q)$ .

The work of [23] defined the notion of max-divergence using  $e$  as the base for logarithm, but in this work we use a variation of it using base 2, which is the same up to a multiplicative constant factor  $\lg e$ . The following lemma lists some of the basic properties of max-divergence (see Definition 2.8).

**Lemma 2.9** (Properties of max-divergence). Let  $X, Y$  be distributions and  $p < 1$ .

1. The following conditions are equivalent.

- (a)  $D_\infty(X||Y) \leq \lg(1/(1-p))$ .
  - (b) For all  $a \in \text{Supp}(X)$  it holds that  $\Pr[X = a] \cdot (1-p) \leq \Pr[Y = a]$ .
  - (c) There exists some random variable  $Z$  such that  $Y = (1-p)X + pZ$ . Namely,  $Y$  can be sampled as: with probability  $1-p$  sample from  $X$  and with probability  $p$  sample from  $Z$ .
2. For  $\text{Supp}(Y) \subseteq \{0, 1\}^m$ ,  $H_\infty(Y) \geq k$  iff  $D_\infty(Y||U_m) \leq m - k$ .
  3. If  $D_\infty(X||Y) \leq r$  and  $D_\infty(Y||X) \leq r$ , then  $D_{\text{KL}}(X||Y) \leq r(2^r - 1)$ .

*Proof Sketch.* Here we only sketch the proofs as they are straightforward. The equivalence of Parts 1a and 1b directly follows from the definition of max-divergence, so here we only show the equivalence of Parts 1b and 1c. Assuming Part 1c we have

$$\Pr[X = a] \cdot (1-p) \leq \Pr[X = a] \cdot (1-p) + \Pr[Z = a] \cdot p = \Pr[Y = a]$$

which implies Part 1b. Assuming Part 1b, we define the distribution  $Z$  over  $\text{Supp}(Y)$  as follows:  $Z(a) = (Y(a) - (1-p) \cdot X(a))/p$ . It is easy to see that  $Z(a) \geq 0$  and that  $\sum_a Z(a) = 1$ , so  $Z$  indeed defines a distribution. Moreover, we have

$$X(a) \cdot (1-p) + Z(a) \cdot p = X(a) \cdot (1-p) + (Y(a) - X(a)) \cdot (1-p) = \Pr[Y = a]$$

which implies that  $Y = (1-p)X + pZ$ , proving Part 1c.

Part 2 directly follows from the definitions of min-entropy and max-divergence.

Part 3 follows from the same proof given in [23] but using the logarithm base 2 instead of  $e$ .  $\square$

## 2.2 Santha-Vazirani Sources and Their Generalizations

**Definition 2.10** (SV sources [39]). A joint distribution  $\bar{X} = (X_1, \dots, X_n)$  where  $X_i \in \{0, 1\}$  for all  $i \in [n]$  is a  $\delta$ -Santha-Vazirani ( $\delta$ -SV) source with bias at most  $\delta \in [0, 1]$ , if for all  $i \in [n]$  and all  $x_1, \dots, x_i \in \{0, 1\}$  it holds that  $(1-\delta)/2 \leq \Pr[X_i = x_i \mid X_1 = x_1, \dots, X_{i-1} = x_{i-1}] \leq (1+\delta)/2$ .

The following definition is a close variant of Block SV Sources defined in [13] where we allow the blocks to have different lengths and specify the amount of *loss* in the min-entropy (compared to the uniform distributing) in each block.

**Definition 2.11** (Block SV Sources [13]). Suppose  $\bar{X} = (X_1, \dots, X_n)$  is a joint distribution where  $X_i \in \{0, 1\}^\ell$  for all  $i \in [n]$ . We call  $\bar{X}$  a  $(\ell, k)$ -block SV source if for all  $i \in [n]$  and all possible  $(x_1, \dots, x_{i-1}) \leftarrow (X_1, \dots, X_{i-1})$  it holds that  $H_\infty(X_i \mid x_1, \dots, x_{i-1}) \geq k$ .

It is easy to see that a  $\delta$ -SV source is a  $(1, 1-\gamma)$ -block-SV source for  $\gamma = \lg(1+\delta) \leq \delta$ . The following definition by Beigi, Etesami and Gohari [6] generalizes the above definitions of SV and Block-SV sources.

**Definition 2.12** (Generalized SV Sources [6]). Let  $\mathcal{D}$  be a set of distributions (dices) over alphabet  $C$ . A distribution  $\bar{X} = (X_1, \dots, X_n)$  over  $C^m$  is a *Generalized SV source* w.r.t  $\mathcal{D}$  if for all  $i \in [n]$  and  $x_1, \dots, x_{i-1} \in C$  there exists  $S \in \text{Span}(\mathcal{D})$  such that for all  $x_i \in C$  it holds that

$$\Pr[X_i = x_i \mid X_1 = x_1, \dots, X_{i-1} = x_{i-1}] = \Pr[S = x_i].$$

### 3 Blockwise $p$ -Tampering: Definitions and Main Results

In this section, we will describe our results formally.

**Notation on sequences of random variables.** By  $D^n$  we denote the product distribution  $D \times \dots \times D$  ( $n$  times). Using this notation, by  $U_m^n$  we mean a sequence of  $n$  blocks each distributed independently like  $U_m$ . Thus, although both of  $U_m^n$  and  $U_n^m$  are eventually  $m \cdot n$  random bits, one is divided into  $n$  blocks and one is divided into  $m$  blocks. For a vector  $x = (x_1, \dots, x_n)$  we let  $x_{\leq i} = (x_1, \dots, x_i)$ ,  $x_{< i} = (x_1, \dots, x_{i-1})$ .

**Definition 3.1** (Valid prefixes and conditional sampling). Let  $\bar{X} = (X_1, \dots, X_n)$  be a joint distribution. We call  $x_{\leq i} = (x_1, \dots, x_i)$  a *valid prefix* for  $\bar{X}$  if there are  $x_{i+1}, \dots, x_n$  such that  $(x_1, \dots, x_n) \in \text{Supp}(\bar{X})$  (i.e.,  $x_{\leq i} \in \text{Supp}(X_{\leq i})$ ). We use  $\text{ValPref}(\bar{X})$  to denote the set of all valid prefixes of  $\bar{X}$  (including the empty string  $x_{\leq 0}$ ). For a valid prefix  $y_{\leq i} \in \text{ValPref}(\bar{X})$ , by  $(X_i \mid y_{\leq i-1})$  we denote the conditional distribution  $(X_i \mid X_1 = y_1, \dots, X_{i-1} = y_{i-1})$ .

**Definition 3.2** (Online-samplable sequences of random variables). We call a *randomized* algorithm  $S(\cdot)$  an *online sampler* for a joint distribution Let  $\bar{X} = (X_1, \dots, X_n)$  if for every valid prefix  $x_{\leq i-1} \in \text{ValPref}(\bar{X})$ , it holds that  $S(x_{\leq i-1})$  outputs according to  $(X_i \mid x_{\leq i-1})$ . If  $\bar{X} = \bar{X}^{(n)}$  is a vector from a *family* of vectors indexed by  $n$ , we let  $N = N(n)$  be the total length of the representation of  $\bar{X}$  (i.e.,  $(X_1, \dots, X_n) \in \{0, 1\}^N$ ) and assume that  $n$  could be derived from  $N(n)$ . In that case, an online sampler  $S(\cdot)$  for  $\bar{X}^{(n)}$  takes also  $N$  as input and it holds that  $S(1^N, x_{\leq i-1}) \equiv (X_i \mid x_{\leq i-1})$ . We call  $\bar{X} = \bar{X}^{(n)}$  *efficiently online samplable* if there exists an online sampler  $S$  for  $\bar{X}$  that runs in polynomial time (i.e.  $\text{poly}(N)$ ). When  $n$  is clear from the context we might simply drop  $1^N$  and simply write  $S(x_{\leq i-1})$ .

**Definition 3.3** (Tampering algorithms for sequences of random variables). Let  $\bar{X} = (X_1, \dots, X_n)$  be an arbitrary joint distribution. We call a (potentially randomized and even computationally unbounded) algorithm  $\text{Tam}$  an (online) *tampering algorithm* for  $\bar{X}$  if given any valid prefix  $x_{\leq i-1} \in \text{ValPref}(\bar{X})$ ,  $\text{Tam}(x_{\leq i-1})$  always outputs  $x_i$  such that  $x_{\leq i} \in \text{ValPref}(\bar{X})$ . If  $\bar{X} = \bar{X}^{(n)}$  is a vector from a *family* of vectors indexed by  $n$ , we call  $\text{Tam}$  an *efficient tampering algorithm* for  $\bar{X}$  if it runs in time  $\text{poly}(N)$  where  $N = N(n)$  is the total bit length of the vector  $\bar{X}$  (i.e.,  $(X_1, \dots, X_n) \in \{0, 1\}^N$ ).

Note that in Definition 3.3, we only allow the tampering algorithm to produce something in the support set of the joint distribution. The following definition defines a notation for representing the “chances” that might be given to a tampering algorithm to tamper with the joint distribution  $\bar{X} = (X_1, \dots, X_n)$ . We need this generalization to formally define the robustness of  $p$ -tampering attack when  $p$  changes during the attack.

**Definition 3.4** (Probability trees over sequences of random variables). Let  $\bar{X} = (X_1, \dots, X_n)$  be an arbitrary joint distribution. We call a function  $\rho: \text{ValPref}(\bar{X}) \mapsto [0, 1]$  a *probability tree* over  $\bar{X}$ . For  $0 \leq p \leq q \leq 1$ , we call  $\rho[\cdot]$  a  $[p, q]$ -probability tree over  $\bar{X}$  if  $\rho(x_{\leq i}) \in [p, q]$  for all  $x_{\leq i} \in \text{ValPref}(\bar{X})$ . We call  $\rho[\cdot]$  the  $p$ -probability tree over  $\bar{X}$  if  $\rho[x_{\leq i}] = p$  for all  $x_{\leq i} \in \text{ValPref}(\bar{X})$ .

Now we define the outcome of an actual “tampering game” in which a tampering algorithm gets to tamper with a joint distribution  $\bar{X} = (X_1, \dots, X_n)$  according to some probability tree defined over  $\bar{X}$ .

**Definition 3.5** ( $\rho$ -tampering variations of distributions). Let  $\bar{X} = (X_1, \dots, X_n)$  be an arbitrary joint distribution, and let  $\rho[\cdot]$  be a probability tree over  $\bar{X}$ . We say that a tampering algorithm  $\text{Tam}$  for  $\bar{X}$  generates  $\bar{Y}$  from  $\bar{X}$  through a  $\rho$ -tampering attack if  $\bar{Y} = (Y_1, \dots, Y_n)$  is inductively sampled as follows. Given any valid prefix  $y_{\leq i-1} \in \text{ValPref}(\bar{Y})$  we will sample  $Y_i$  through the following process:

- with probability  $1 - \rho[y_{\leq i-1}]$ , sample  $Y_i$  from  $(X_i \mid X_{\leq i-1} = y_{\leq i-1})$ , and
- with probability  $\rho[y_{\leq i-1}]$ , sample  $Y_i \leftarrow \text{Tam}(y_{\leq i-1})$ .

Equivalently, using Definition 2.9, for all  $y_{\leq i-1} \in \text{ValPref}(\bar{Y})$  we have  $(Y_i \mid y_{\leq i-1}) = (1 - \rho[y_{\leq i-1}]) \cdot (X_i \mid X_{\leq i-1} = y_{\leq i-1}) + \rho[y_{\leq i-1}] \cdot \text{Tam}(y_{\leq i-1})$ . In this case, we also call  $\bar{Y}$  a  $\rho$ -tampering variation of  $\bar{X}$ . In case  $\rho$  is the constant function  $p$ , we call  $\bar{Y}$  a  $p$ -tampering variation of  $\bar{X}$  and we say that Tam generates  $\bar{Y}$  from  $\bar{X}$  through a  $p$ -tampering attack.

Note that even in cases where we end up sampling  $Y_i$  from the “untampered” distribution of  $X_i$  (which happens with probability at least  $1 - \rho[x_{\leq i-1}]$ ) we still sample from  $X_i$  conditioned on the *possibly tampered* prefix  $(y_1, \dots, y_i)$ . In other words, the result of the tampering algorithm determines, in case it happens, will completely substitute the tampered block and the sampling will continue as if the history of the blocks were from the untampered sequence  $X_1, \dots, X_i$ . For the special case that  $X_i$ 's are independent distributions (e.g., when  $\bar{X}$  is uniform distribution over some set  $\Sigma^n$ ) we will not need to do this.

**Prefixes remain valid.** Note that because in Definition 3.5 the algorithm Tam is a (valid) tampering algorithm for  $\bar{X}$ , all the resulting prefixes will remain valid for  $\bar{X}$  and we will have  $\text{ValPref}(\bar{Y}) \subseteq \text{ValPref}(\bar{X})$ . In fact, we get  $\text{ValPref}(\bar{Y}) = \text{ValPref}(\bar{X})$  if  $\rho[x_{\leq i}] < 1$  for all  $x_{\leq i} \in \text{ValPref}(\bar{X})$ . A more general definition of tampering algorithms (compared to Definition 3.3) could use a larger support set  $\mathcal{Z}$  where  $\text{ValPref}(\bar{X}) \subset \mathcal{Z}$  and only require the tampering algorithm to produce prefixes in  $\mathcal{Z}$ . However, since our main contributions in this paper is to give attacks, by restricting our model to require the attackers to remain in  $\text{ValPref}(\bar{X})$  only makes our results stronger.

**Remark 3.6** (Efficient tampering vs. efficient sampling). Note that an *efficient tampering* refers only to when the algorithm Tam is polynomial time, and it can apply even to settings where  $\bar{X}$  and its variation generated by Tam are *not* efficiently samplable. On the other hand, using the standard terminology,  $\bar{X}$  is efficiently samplable if one can efficiently sample *all* of the blocks of  $\bar{X}$  *simultaneously*. Of course, if  $\bar{X}$  is efficiently *online* samplable and if Tam is also an efficient tampering for  $\bar{X}$ , then the variation  $\bar{Y}$  of  $\bar{X}$  produced by tampering attack Tam will also be trivially efficiently online-samplable, but we emphasize that this is a specific way of getting an efficient sampler for  $\bar{Y}$ , and so the efficiency of our tampering attacks shall not be confused with mere efficient samplability of the final distribution  $\bar{Y}$ .

**Remark 3.7.** An alternative variant of Definition 3.5 could ‘strengthen’ the tampering algorithm Tam who, now, receives the ‘original’ sample  $x_i$  before substituting it with something else. Namely, we would first sample  $x_i \leftarrow (X_i \mid y_{\leq i-1})$ , and then with probability  $1 - p$  we let  $y_i = x_i$  and with probability  $p$  we let  $y_i = \text{Tam}(y_{\leq i-1}, x_i)$ . This definition is natural for scenarios in which the adversary gets to see the first initial sample and then might decide to change or not change it. However, as long as either (1) tampering is allowed to be inefficient or (2)  $\bar{X}$  is efficiently online samplable, the power of tampering attacks under this alternative definition is the same as those under Definition 3.5. To see why, first note that  $\text{Tam}(y_{\leq i-1}, x_i)$  can always ignore the extra input  $x_i$ . In the other direction, suppose  $\text{Tam}'$  is a tampering algorithm under the alternative definition and suppose a tampering algorithm  $\text{Tam}(y_{\leq i-1})$  is only given  $y_{\leq i-1}$ . If Tam can obtain a sample  $x'_i \leftarrow (X_i \mid y_{\leq i-1})$ , then it could also emulate  $\text{Tam}'(y_{\leq i-1}, x'_i)$ . Interestingly, although  $x_i$  and  $x'_i$  might be different samples, this emulation of  $\text{Tam}'(y_{\leq i-1}, x'_i)$  by Tam leads to the same distribution.

Now we define what it means for a tampering adversary to successfully bias the output of a function, while being robust to changes in probabilities.

**Definition 3.8** (Robust  $p$ -tampering attacks for biasing real functions). Let  $\bar{X} = (X_1, \dots, X_n)$  be a joint distribution,  $f: \text{Supp}(\bar{X}) \mapsto \mathbb{R}$  a real function and Tam a tampering algorithm for  $\bar{X}$ .

- For a probability tree  $\rho$  over  $\bar{X}$ , we say that Tam is a  $\rho$ -tampering attack biasing  $f(\bar{X})$  by at least  $\delta$ , if Tam generates  $\bar{Y}$  from  $\bar{X}$  through a  $\rho$ -tampering attack and  $\mathbb{E}[f(\bar{Y})] \geq \mathbb{E}[f(\bar{X})] + \delta$ .
- For  $p \in [0, 1]$ , we say that Tam is a  $p$ -tampering attack biasing  $f(\bar{X})$  by at least  $\delta$ , if Tam a  $\rho$ -tampering attack biasing  $f(\bar{X})$  by at least  $\delta$  for the constant probability tree  $\rho[x_{\leq i}] = p$ .
- We say that Tam is a *robust*  $p$ -tampering attack biasing  $f(\bar{X})$  by at least  $\delta$ , if for *every*  $[p, 1]$ -probability tree  $\rho$  over  $\bar{X}$  it holds that Tam is a  $\rho$ -tampering attack biasing  $f(\bar{X})$  by at least  $\delta$ .

### 3.1 Main Results: Blockwise $p$ -Tampering of Bounded Functions

Now, we are ready our main results that are about biasing real functions through *efficient* blockwise  $p$ -tampering attacks. We will then describe our results about the computationally unbounded setting where the tampering algorithm Tam is not necessarily polynomial time. Our main motivation for studying the computationally unbounded setting is to understand the *limitations* of what amount of bias could be achieved. We will then describe the applications of our results for attacking candidate randomness extractors (over multiple sources or variations of SV sources) through  $p$  tampering attacks.

**Theorem 3.9** (Efficient blockwise  $p$ -tampering of bounded real functions). *Let  $\bar{X} = (X_1, \dots, X_n)$  be a joint distribution,  $f: \text{Supp}(\bar{X}) \mapsto [-1, +1]$  be a real-output function defined over  $\text{Supp}(\bar{X})$ . Then there is a tampering algorithm Tam for  $\bar{X}$  such that:*

1. **(Bias)** Tam is a robust  $p$ -tampering attack biasing  $f(\bar{X})$  by at least  $\Omega(p) \cdot \text{Var}[f(\bar{X})]$ . In particular, if the function  $f: \text{Supp}(\bar{X}) \mapsto \{-1, +1\}$  is Boolean, then the bias is at least  $\frac{p}{2+2p} \cdot \text{Var}[f(\bar{X})]$ .
2. **(Efficiency)** Moreover, Tam could be implemented efficiently given oracle access to any online sampler  $S(\cdot)$  for  $\bar{X}$  and  $f(\cdot)$ . In particular, given only two samples  $y_i^1, y_i^2 \leftarrow S(y_{\leq i-1})$ , Tam( $y_{\leq i-1}$ ) chooses between  $y_i^1, y_i^2$  by making use of a biased coin that only depends on  $f[y_{\leq i-1}, y_i^1]$ . Such biased coin could be sampled efficiently using further calls to  $S(\cdot)$  and one call to  $f(\cdot)$ .

Theorem 3.9 above extends the previous result of [1] from bitwise to blockwise  $p$ -tampering, though with worse constants. See Subsection 5.1 for the full proof of Theorem 3.9.

**Importance of the efficiency features of the attacker in Theorem 3.9.** As we will see in Theorem 3.11 below, we can get better biasing bounds for the Boolean case than  $p \cdot \text{Var}[f(\bar{X})]/4$ , however, the reason that we pointed this out in Theorem 3.9 was that result comes along with the efficiency feature specified in Theorem 3.9 (and this is not the case in our Theorem 3.11 below). As mentioned, the attacker of Theorem 3.9 only needs *two honestly* generated samples  $\{y_i^1, y_i^2\}$  for the next tampered block  $X_i$  and chooses one of them. Interestingly, this means that if the tampering algorithm is actually given an ‘initial true value’  $x_i$  for block  $X_i$  (e.g., the honestly generated randomness to be used in a randomized algorithm) then the tampering algorithm could basically just either keep  $x_i$  or substitute it with another fresh sample from  $X_i$ . This is a natural attack strategy when the adversary can “reset” the sampling procedure for the block  $X_i$ .

**Biasing Martingales.** An interesting special case of Theorem 3.9 is when the joint distribution  $\bar{X} = (X_1, \dots, X_n)$  is a martingale (i.e.,  $X_i \in \mathbb{R}$  and  $\mathbb{E}[X_i | x_{\leq i-1}] = x_{i-1}$ ) and  $f(\bar{X}) = X_n \in [-1, +1]$ . In this case, it holds that  $\hat{f}[x_{\leq i}] = x_i$ , and so our attacker of Theorem 3.9 becomes extremely simple: given any two samples  $y_i^1, y_i^2 \leftarrow (X_i | y_{\leq i-1})$ ,  $\text{Tam}(y_{\leq i-1})$  chooses  $y_i = y_i^1$  with a probability that only depends on  $y_i^1$  and chooses  $y_i = y_i^2$  otherwise. Note that *no* further calls to the online sampler nor  $f(\cdot)$  is needed! Moreover, this simple attack not only biases the final value  $X_n = f(\bar{X})$  but it does bias *every* other  $X_i$  as well. The reason is that if we define  $f_i(X_{\leq i}) = X_i \in [-1, +1]$ , then the attacker’s algorithm would be identical for biasing  $f_i(\cdot)$  compared to biasing  $f_n(\cdot) = f(\cdot)$ . Therefore, our attack generates a  $p$ -tampering variation  $\bar{Y}$  of  $\bar{X}$  that *simultaneously* achieves bias  $Y_i \geq X_i + \Omega(p) \cdot \text{Var}[X_i]$  for *every* block  $i \in [n]$ . Moreover, the  $p$ -tampering is efficient if the martingale is online samplable.

**Tampering with only a part of randomness.** The specific way that the attacker of Theorem 3.9 chooses between the two samples  $\{y_i^1, y_i^2\}$  for block  $X_i$  allows us to generalize the attack to settings where the tamping happens only over *part* of the randomness and some subsequent randomness  $R$  is also used for computing  $f$ . As we will see, this corollary would also be useful for attacking randomized learners through the so called ‘targeted poisoning’ attacks.

**Corollary 3.10** (Biasing bounded ‘randomized’ functions). *Let  $\bar{X} = (X_1, \dots, X_n)$  be a joint distribution,  $R$  another distribution, and  $f: \text{Supp}(\bar{X} \times R) \mapsto [-1, +1]$ . For any fixed  $x \leftarrow \bar{X}$ , let  $g(x) = \mathbb{E}_{r \leftarrow R}[f(x, r)] \in [-1, +1]$ . Then there is a tampering algorithm  $\text{Tam}$  for  $\bar{X}$  (not receiving  $R$ ) such that:*

1. **(Bias)**  $\text{Tam}$  is a robust  $p$ -tampering attack biasing  $g(\bar{X})$  by at least  $\Omega(p) \cdot \text{Var}[g(\bar{X})]$ .
2. **(Efficiency)**  $\text{Tam}$  could be implemented efficiently given oracle access to any online sampler  $S(\cdot)$  for  $\bar{X}$  and  $f(\cdot, \cdot)$ . In particular,  $\text{Tam}(y_{\leq i-1})$  again chooses between two samples  $y_i^1, y_i^2 \leftarrow S(y_{\leq i-1})$  using further calls to  $S(\cdot)$  and one call to  $f(\cdot, \cdot)$  and one sample from  $R$ .

*Proof of Corollary 3.10 using Theorem 3.9.* To derive Corollary 3.10 from Theorem 3.9 we apply Theorem 3.9 directly to the function  $g(x) = \mathbb{E}[f(x, R)]$ , and we rely on the properties specified in the efficiency part of Theorem 3.9 to derive the efficiency of the new attacker. All we need is to provide a sample from the distribution  $Z$  (for choosing between  $y_i^1, y_i^2 \leftarrow S(y_{\leq i-1})$ ) when we try to bias  $g$ . In order to do so, we can first sample  $x \leftarrow (\bar{X} | y_{\leq i-1}, y_i^1)$  using  $S(\cdot)$ , and then output  $Z \leftarrow f(x, R)$  using one sample  $r \leftarrow R$ . By the linearity of expectation, even though we did not really compute  $g(x)$ , this way of sampling  $Z$  using only one  $r \leftarrow R$  has the needed properties for the (average) function  $g$  as well.  $\square$

The following theorem gives a better biasing bound for the important special case of Boolean functions. On the down side, the attacker will be less efficient and asks more queries to the online sampler  $S(\cdot)$ .<sup>5</sup>

**Theorem 3.11** (Biasing Attacks on Boolean functions). *Let  $\bar{X} = (X_1, \dots, X_n)$  be a joint distribution,  $f: \text{Supp}(\bar{X}) \mapsto \{+1, -1\}$  a Boolean function defined over  $\text{Supp}(\bar{X})$ , and  $\mu = \mathbb{E}[f(\bar{X})]$ . Suppose  $S$  is a sampler for  $\bar{X}$  and let  $N$  be an upper bound on the total binary length of  $\bar{X} = (X_1, \dots, X_n) \in \{0, 1\}^N$ , and  $\varepsilon < 1$  be an input parameter. Then there is a tampering algorithm  $\text{Tam}$  for  $\bar{X}$  that:*

1. **(Bias)**  $\text{Tam}$  is a robust  $p$ -tampering attack biasing  $f(\bar{X})$  by  $\geq \frac{p \cdot (1 - \mu^2)}{2 - p \cdot (1 - \mu)} - \frac{\varepsilon}{1 + \mu}$ .<sup>6</sup>

<sup>5</sup>The sample complexity measure is an important factor in some of the applications of our biasing attacks. For example, to attack the soundness of learning algorithms through targeted poisoning attacks, the sample complexity of the attacker translates into how much ‘fresh’ data is needed to substitute the original training examples when the tampering happens.

<sup>6</sup>The analysis of the greedy attack of [1] shows that the amount of bias is at least  $p \cdot (1 - |\mu|)$ . Our bound depends on  $1 - \mu^2$  instead of  $1 - |\mu|$ . The reason behind this is that we use a better approximation of the probabilities for the output to be  $-1$  or  $+1$ .

2. **(Efficiency)** Moreover, Tam could be implemented in time  $\text{poly}(N/\varepsilon)$  given oracle access to any online sampler  $S(\cdot)$  for  $\bar{X}$  and  $f(\cdot)$ . Thus, if  $\varepsilon \geq 1/\text{poly}(N)$ ,  $\bar{X}$  is efficiently online samplable, and  $f$  is efficient, then Tam would be efficient as well.

We prove our Theorem 3.11 using ideas from the attack of [1] also for the Boolean case. In a nutshell, we follow the same ‘greedy’ approach, but the analysis of the attack in the blockwise setting becomes more challenging and we can no longer get the same bias of  $+p$  in the balanced case. Indeed, as we will show in Section C, achieving the bias of  $+p$  for balanced functions in the blockwise setting is *not* possible in general! See Section 5 (in particular Section 5.2) for the full proof of Theorem 3.11.

**Remark 3.12** (Robustness vs.  $p$ -obliviousness). Note that in both Theorems 3.11 and 3.9 the attackers are robust in the sense that they work simultaneously for all  $[p, 1]$  probability trees (i.e., they only rely on the lower-bound  $p$  for the probability of the tampering to happen for each block). However, this feature of the attacker should not be confused with another aspect of our attackers that they are  $p$ -oblivious, meaning the tampering algorithm Tam does not rely on knowing  $p$  either. Putting these two together, it means that the attackers of Theorems 3.11 and 3.9 could be “generated” independently of the probability tree  $\rho$  under which the tampering to the randomness will eventually happen, and yet the quality of obtained bias only depend on the minimum over all the probabilities under which the blocks become tamperable.

**Computationally Unbounded  $p$ -Tampering.** One might wonder what power of blockwise  $p$ -tampering attacks. Even though our focus in this work is on the computationally bounded setting, we also study the power and limitations of computationally unbounded  $p$ -tampering attacks. Showing the power of attackers in the unbounded model might eventually shed light into how to get better efficient attackers as well, and proving limitations in this model imply strong limits for efficient tampering algorithms as well. In Section C we show that the better biasing bound of Theorem 3.11 could be obtained for bounded real functions as well, but this comes with an inefficient  $p$ -tampering, and achieving this bound efficiently remains as an open question. Perhaps surprisingly, we also show that there are balanced functions over block sources where the best bias by (even inefficient)  $p$ -tampering attacks is smaller than  $0.7p$ . This comes in contrast with the bitwise  $p$ -tampering model where  $p$  is the optimal possible bias in general. See Section C for more details.

## 4 Applications of $p$ -Tampering Biasing Attacks

In this subsection we describe some of the applications of our main results on blockwise  $p$ -tampering of bounded functions in several different contexts.

### 4.1 Efficient $p$ -Tampering Attacks on Extractors

Rather than proving Theorem 1.3, here we prove a more general result by defining yet another generalization of SV sources based on the notion of max-divergence [23] (see Definition 2.8) which is tightly related to  $p$ -tampering variations. Intuitively, we will show that  $\bar{X}$  is an  $(\ell, \gamma)$  block SV source if the uniform distribution  $U_\ell^n$  is a  $p$ -tampering variation of  $\bar{X}$  for  $p \approx \gamma$ . We will then show that our  $p$ -tampering attacker of Theorem 3.9 produces  $\bar{Y}$  such that  $\bar{X}$  itself is a  $O(p)$ -tampering variation of  $\bar{Y}$ . We first define the following generalization of block-SV sources based on max-divergence.

**Definition 4.1** (MD and MMD Sources). Let  $\bar{X} = (X_1, \dots, X_n)$  be a joint distribution. For real number  $r \geq 0$ , we call a joint distribution  $\bar{Y} = (Y_1, \dots, Y_n)$  an  $(\bar{X}, r)$ -max-divergence (MD) source if  $\text{Supp}(\bar{Y}) =$



$\text{Supp}(\bar{X})$  and for all  $i \in [n], x_{<i} \in \text{ValPref}(\bar{X})$  the max-divergence  $D_\infty((X_i | x_{<i}) || (Y_i | x_{<i}))$  is at most  $r$ . We call  $\bar{Y}$  an  $(\bar{X}, r)$  *mutual MD (MMD) source* if in addition  $\bar{X}$  is an  $(\bar{Y}, r)$  MD source as well.

**Remark 4.2** (Sources based on other distance measures). The above definition uses max-divergence in order to limit how ‘far’ the source  $\bar{Y}$  can be from the ‘central’ random process  $\bar{X} = (X_1, \dots, X_n)$ . Alternative definitions could be obtained by using other distance metrics and measures. For example, we can also define  $(\bar{X}, r)$  KL sources to include all distributions  $\bar{Y}$  such that  $D_{\text{KL}}((X_i | x_{<i}) || (Y_i | x_{<i})) \leq r$ . A result of [23] (see Part 3 of Lemma 2.9) shows that any  $(\bar{X}, r)$  *mutual MD source* is also a  $(\bar{X}, r')$  KL-source for  $r' = r(2^r - 1)$  which is  $r' \leq r^2$  for any  $r \leq 1$ .

The following claim shows that MD sources and  $p$ -tampering variations are tightly related. The proof directly follows from definitions of MD sources and  $p$ -variations.

**Claim 1** (MD sources vs. tampering variations).  $\bar{Y} = (Y_1, \dots, Y_n)$  is an  $(\bar{X}, r)$ -MD source iff it is a  $p$ -tampering variation of  $\bar{X}$  for  $p = 1 - 2^{-r}$ .

The following claim shows that MD sources are also related to SV block sources (in the ‘reverse’ direction), and its proof directly follows from the definition of MD sources and Part 2 of Lemma 2.9.

**Claim 2** (MD sources vs. block SV sources). For a joint distribution  $\bar{X} = (X_1, \dots, X_n)$ ,  $U_\ell^n$  is an  $(\bar{X}, r)$ -MD source iff  $\bar{X}$  is an  $(\ell, \ell - r)$  block SV source. In particular, if  $\bar{X}$  is an  $(U_\ell^n, \ell - r)$ -MMD source, then it is also an  $(\ell, \ell - r)$ -block SV source.

Theorem 1.3 follows from Claim 2 above and the following general result about the impossibility of deterministic extraction from MMD sources.

**Theorem 4.3** (Impossibility of extractors from MMD sources). Let  $\bar{X} = (X_1, \dots, X_n)$  be a joint distribution with an efficient online sampler, and let  $f: \text{Supp}(\bar{X}) \mapsto \{+1, -1\}$  be an efficient Boolean function. Then, there is a  $p$ -tampering variation  $\bar{Y}$  of  $\bar{X}$  where:

1.  $\bar{Y}$  is an  $(\bar{X}, p)$  MMD source.
2.  $|\mathbb{E}[f(\bar{Y})]| \geq \Omega(p)$ .
3.  $\bar{Y}$  is generated by an efficient tampering algorithm Tam.

The first two items in Theorem 4.3 imply that  $f$  cannot be an extractor for  $(\bar{X}, p)$  MMD sources for any  $\bar{X} = (X_1, \dots, X_n)$ . Moreover, one can show that the source  $\bar{Y}$  is also a  $(\bar{X}, p^2)$  KL source because it is a  $(\bar{X}, p)$  *mutual MD source* (see Remark 4.2).

**Efficiency of the attacker.** The last condition shows that the  $p$ -tampering attack against such  $f$  (as a candidate extractor) could be implemented by an *efficient*  $p$ -tampering attacker. We emphasize that the efficiency condition again is crucial here. In fact, if we change the statement of Theorem 4.3 by (1) restricting  $\bar{X} = (Z \times \dots \times Z)$  to iid distributions and more importantly (2) allowing Tam to be computationally unbounded, then we can derive this weaker version of Theorem 4.3 from the recent impossibility result of [6] for generalized SV sources as follows. Beigi et al. [6] showed that bit extraction with  $o(1)$  bias from generalized SV sources (Definition 2.12) is impossible if (1) all the distributions  $D \in \mathcal{D}$  available to the adversary have full support over the alphabet set  $C$  and that (2) the span of distributions  $\mathcal{D}$  (see Definition 2.3) has full dimension  $|C|$ . To apply their result to MMD sources, we observe that (1) the distribution of  $Y_i$  where  $\bar{Y} = (Y_1, \dots, Y_n)$  is an  $(\bar{X}, r)$  MMD source has full support (i.e.,  $\text{Supp}(Z) = C$ ) and that (2) conditioned on any  $y_{\leq i-1}$ , the set of all possible distributions for  $Y_i$  forms a polytope with rank  $|\text{Supp}(Z)|$ .

*Proof of Theorem 4.3.* To prove Theorem 4.3 we use Theorem 3.9 and rely on some specific properties of the  $p$ -tampering attacker there. Even though the function  $f$  is Boolean, for some minor technical reasons, we will actually use the  $p$ -tampering attacker of Theorem 3.9 for *real* output functions, which is a bit different. In the following we will show that this attacker has the properties listed in Theorem 4.3.

First note that without loss of generality, we can assume that  $\mathbb{E}[f(\bar{X})] \geq 0$  (as otherwise we can work with  $-f$  and bias it towards  $+1$ ). In that case, the second and third properties of Theorem 4.3 follow from the main properties of Tam as stated in Theorem 3.9. However, for getting the first property (that it gives us an MMD source) we need to get into the actual attack's description from the proof of Theorem 3.9 given in Subsection 5.1, which we also describe here. This attacker Tam (for the *real* output case) is based on one-rejection sampling (of Construction 5.6) modified as follows. Whenever the tampering algorithm is given the chance to tamper with a new block (which happens with probability  $p$ ), the attacker itself tosses a coin and decides *not* to tamper with the block with probability 0.5, and otherwise will actually run the one-rejection sampling of Construction 5.6. Thus, during the execution of the  $p$ -tampering attack, the tampering actually happens with probability  $p/2$ .

As described above, the tampering happens with probability  $p/2$ , so by Claim 1, it holds that  $\bar{Y}$  is an  $(\bar{X}, r)$  MD source for  $r \leq \lg(1/(1 - p/2)) \leq p$  (by  $p \in [0, 1]$ ). On the other hand, the one-rejection sampling is actually used only with probability  $p/2$ . Therefore, for every possible  $y_{\leq i}$ , if we let  $\alpha = \Pr[X_i = y_i \mid y_{\leq i-1}]$ , then it holds that  $\Pr[Y_i = y_i \mid y_{\leq i-1}] \leq (1 - p/2) \cdot \alpha + (p/2) \cdot (2\alpha) \leq (1 + p/2) \cdot \alpha$ , because, either no tampering happens with probability  $1 - p/2$  and even if it happens, because the tampering algorithm only uses two samples for the tampered block, by a union bound, the probability of sampling  $y_i$  in this case is at most  $2\alpha$ , which means that  $\bar{X}$  is an  $(\bar{Y}, r)$  MD source for  $r \leq \lg(1 + p/2) \leq p$  (by  $p \in [0, 1]$ ).

Putting things together, it holds that  $\bar{Y}$  is indeed an  $(\bar{X}, p)$  MMD source.  $\square$

## 4.2 Targeted Poisoning Attacks on Learners

**Terminology.** Let  $\mathcal{D}$  be the domain containing all the objects of interest in a learning problem, and let  $\mathcal{C}$  be a class of *concept* functions mapping objects in  $\mathcal{D}$  to a set of labels  $\mathcal{T}$ . A labeled example from the set  $\mathcal{D}$  for a concept function  $c \in \mathcal{C}$  is a pair  $x = (d, c(d))$  where  $d \in \mathcal{D}$ . We use  $\mathcal{P}_c = \{(d, c(d)) \mid d \in \mathcal{D}\}$  to denote all the labeled examples from  $\mathcal{D}$ . The goal of a learning algorithm  $L$  is to produce a *hypothesis*  $h \in \mathcal{H}$  after receiving a sequence  $x = (x_1, \dots, x_n)$  of labeled examples that we call the training sequence, such that  $h$  can predict the label of a given input from  $\mathcal{D}$ . The examples in the training sequence are usually sampled independently from a distribution  $X$  over  $\mathcal{P}_c$  through an oracle  $O_X(\cdot)$  that we call the *training oracle*. A subset  $\mathcal{X} \subseteq \mathcal{P}_c$  is a *test set* if we use it to evaluate the performance of the hypothesis  $h$ .

**Definition 4.4** (Cost and average cost). A *cost function*  $\text{cost} : \mathcal{H} \times 2^{\mathcal{P}_c} \rightarrow [0, 1]$  captures the quality of a hypothesis, and the lower the value of  $\text{cost}(h, \mathcal{X})$ , the better  $h$  is performing on the examples in  $\mathcal{X}$ . We define the average cost function for a learning algorithm  $L$  and a test set  $\mathcal{X}$  according to a specific training oracle as follows:

$$\overline{\text{cost}}_L^O(\mathcal{X}) = \mathbb{E}_{\substack{x_1, \dots, x_n \leftarrow O \\ h \leftarrow L(x_1, \dots, x_n)}} [\text{cost}(h, \mathcal{X})]$$

For example the cost functions might be the fraction of examples in  $\mathcal{X}$  that  $h$  generate a wrong label for. The test set itself can consist of only one point, or it might be very large to model the scenario where sampling an example from  $\mathcal{X}$  is equivalent to sampling from  $X$ .<sup>7</sup>

<sup>7</sup>In case the test data comes from  $X$  itself (i.e.,  $\mathcal{X} \equiv X$ ), the average cost becomes tightly related to PAC learnability [42]. In particular, if we define cost to be one whenever the hypothesis  $h$  generates a wrong label, then any  $(\varepsilon, \delta)$ -PAC learner has average cost at most  $\varepsilon + \delta$ . Conversely, if the average cost is at most  $\gamma$ , then by an averaging argument we get a  $(\sqrt{\gamma}, \sqrt{\gamma})$ -PAC learner.

**Definition 4.5** ( $p$ -tampering training oracles). Let  $O_X$  be the training oracle for a distribution  $X$ . A  $p$ -tampering oracle  $\widehat{O}_X^p$  works as follows. Whenever the training algorithm queries this oracle, with probability  $1-p$  the answer is generated from the original oracle  $O_X$  and with probability  $p$  a stateful adversary gets the control over the oracle and answers with an arbitrary pair  $(d, t)$  such that  $(d, t) \in \mathcal{P}_c$ . We call  $\widehat{O}_X^p$  efficient, if the pair  $(d, t)$  is generated using an efficient  $p$ -tampering algorithm that takes as input  $1^N$ , where  $N$  is the total length of the training sequence  $x$ , and all the previous samples in the training sequence.

We can use our Theorem 3.9 to increase the average cost of even randomized learners where the cost could also be a real number. In the following theorem we do exactly that. However, the quality of this attack depends on the variance of the learners success probability (as defined in Theorem 4.6). Thus, a provably secure randomization defense against our attacks need, as the very least, to upper bound the variance parameter defined in Theorem 4.6.

**Theorem 4.6** (Power of targeted poisoning attack against real cost functions). *Let  $\mathcal{C}$  be a concept class defined over domain  $\mathcal{D}$ . Also let  $L$  be a (potentially randomized) learning algorithm for  $\mathcal{C}$  which takes a sequence of labeled examples  $x = (x_1, \dots, x_n)$  that are sampled using an efficient training oracle  $O_X$  and outputs a hypothesis  $h \in \mathcal{H}$ . For any such learning algorithm  $L$  that tries to learn a concept  $c \in \mathcal{C}$ , any  $p \in [0, 1]$ , any test set  $\mathcal{X}$  and any cost function  $\text{cost} : \mathcal{H} \times 2^{\mathcal{P}_c} \rightarrow [0, 1]$  there exists a  $p$ -tampering training oracle  $\widehat{O}_X^p$  such that if we sample  $x$  using  $\widehat{O}_X^p$  instead of  $O_X$  the average cost increases as*

$$\overline{\text{cost}}_{L^{\widehat{O}_X^p}}(\mathcal{X}) \geq \overline{\text{cost}}_{L^{O_X}}(\mathcal{X}) + \Omega(p \cdot \sigma^2) \text{ where } \sigma^2 = \text{Var}_{x_1, \dots, x_n \leftarrow O_X} \left[ \mathbb{E}_{h \leftarrow L(x_1, \dots, x_n)} [\text{cost}(h, \mathcal{X})] \right].$$

Moreover, if  $L$  is efficient,  $X$  is efficiently samplable, and  $\text{cost}(\cdot)$  is efficiently computable, then the corresponding  $p$ -tampering attack is efficient as well.

*Proof.* Assume  $L$  uses its own randomness  $r \leftarrow R$  in addition to  $(x_1, \dots, x_n)$  and outputs a hypothesis  $h$ . For a fixed test set  $\mathcal{X}$ , we define a function  $f : \mathcal{C}_p^n \times \text{Supp}(R) \rightarrow [-1, +1]$  as follows:

$$f(x_1, \dots, x_n, r) = 2 \cdot \text{cost}(L(x_1, \dots, x_n, r), \mathcal{X}) - 1.$$

The output of the cost function is between 0 and 1, so the output of  $f$  is between  $-1$  and  $+1$ . Now by using our biasing attacks over *part* of the randomness of randomized functions (i.e., Corollary 3.10) there exists a  $p$ -tampering variation  $\overline{Y}$  of  $X^n$ , generated through an efficient tampering attack, that biases  $f$  as follows:

$$\widehat{\mu} = \mathbb{E}_{\substack{x_1, \dots, x_n \leftarrow \overline{Y} \\ r \leftarrow R}} [f(x_1, \dots, x_n, r)] > \mu + \Omega(p) \cdot v$$

$$\text{where } \mu = \mathbb{E}_{\substack{x_1, \dots, x_n \leftarrow X^n \\ r \leftarrow R}} [f(x_1, \dots, x_n, r)] \text{ and } v = \text{Var}_{x_1, \dots, x_n \leftarrow X^n} \left[ \mathbb{E}_{r \leftarrow R} [f(x_1, \dots, x_n, r)] \right].$$

Since  $\overline{Y}$  is a  $p$ -tampering variation of  $X^n$  generated by an efficient tampering attack, there is an efficient  $p$ -tampering training oracle  $\widehat{O}_X^p$  that generates  $\overline{Y}$ . By the linearity of expectation,  $\widehat{\mu} = 2 \cdot \overline{\text{cost}}_{L^{\widehat{O}_X^p}}(\mathcal{X}) - 1$ ,  $\mu = 2 \cdot \overline{\text{cost}}_{L^{O_X}}(\mathcal{X}) - 1$ . In addition, we have  $v = 4 \cdot \sigma^2$ , so by replacing  $\widehat{\mu}$ ,  $\mu$  and  $v$  we get

$$\overline{\text{cost}}_{L^{\widehat{O}_X^p}}(\mathcal{X}) \geq \overline{\text{cost}}_{L^{O_X}}(\mathcal{X}) + \Omega(p) \cdot \sigma^2.$$

□

This bound of the above theorem could be indeed very weak as it depends on the variance of the cost of the generated hypothesis. In particular, the change could be  $o(1)$ . As we will see, for the special case of Boolean cost functions (e.g., classification) we can increase the error arbitrarily close to one.

**Theorem 4.7** (Power of targeted poisoning attacks against classifiers). *Let  $\mathcal{C}$  be a concept class defined over domain  $\mathcal{D}$ . Also let  $L$  be a deterministic, learning algorithm for  $\mathcal{C}$  which takes a sequence of labeled examples  $x = (x_1, \dots, x_n)$  that are sampled using an efficient training oracle  $O_X$  and outputs a hypothesis  $h \in \mathcal{H}$ . For any such learning algorithm  $L$  that tries to learn a concept  $c \in \mathcal{C}$ , any  $p \in [0, 1]$ , any test set  $\mathcal{X}$  and any cost function  $\text{cost} : \mathcal{H} \times 2^{\mathcal{P}^c} \rightarrow \{0, 1\}$  there exist a  $p$ -tampering training oracle  $\widehat{O}_X^p$  such that if we sample  $x$  using  $\widehat{O}_X^p$  instead of  $O_X$ , the average cost increases as:*

$$\overline{\text{cost}}_{L^{\widehat{O}_X^p}}(\mathcal{X}) \geq \delta + \frac{p \cdot (\delta - \delta^2)}{1 - p \cdot (1 - \delta)} \quad \text{where} \quad \delta = \overline{\text{cost}}_{L^{O_X}}(\mathcal{X}).$$

Moreover, if  $L$  and  $\text{cost}(\cdot)$  are efficient and  $X$  is efficiently samplable, then for any  $\varepsilon > 0$  our  $p$ -tampering training oracle can be implemented in time  $\text{poly}(\frac{n}{\varepsilon \delta})$  and achieve  $\overline{\text{cost}}_{L^{\widehat{O}_X^p}}(\mathcal{X}) \geq \delta + \frac{p \cdot (\delta - \delta^2)}{1 - p \cdot (1 - \delta)} - \varepsilon$ .

The proof of Theorem 4.7 is based on Theorem 3.11.

*Proof of Theorem 4.7.* We define a function  $f : \mathcal{C}^n \rightarrow [-1, +1]$  as follows:

$$f(x_1, \dots, x_n) = 2 \cdot \text{cost}(L(x_1, \dots, x_n), \mathcal{X}) - 1.$$

Now using Theorem 3.11, there exist a  $p$ -tampering variation  $\overline{Y}$  of  $X^n$  that biases  $f$  as follows:

$$\widehat{\mu} = \mathbb{E}_{x_1, \dots, x_n \leftarrow \overline{Y}} \geq \mu + \frac{p \cdot (1 - \mu^2)}{2 - p \cdot (1 - \mu)} \quad \text{where} \quad \mu = \mathbb{E}_{x_1, \dots, x_n \leftarrow X^n} [f(x_1, \dots, x_n)].$$

Since  $\overline{Y}$  is a  $p$ -tampering variation of  $X^n$ , there is an  $p$ -tampering training oracle  $\widehat{O}_X^p$  that generates  $\overline{Y}$ . With a simple calculation we have  $\widehat{\mu} = 2 \cdot \overline{\text{cost}}_{L^{\widehat{O}_X^p}}(\mathcal{X}) - 1$  and  $\mu = 2 \cdot \delta - 1$ . By replacing  $\widehat{\mu}$  and  $\mu$  we get

$$\overline{\text{cost}}_{L^{\widehat{O}_X^p}}(\mathcal{X}) \geq \delta + \frac{p \cdot (\delta - \delta^2)}{1 - p \cdot (1 - \delta)}.$$

The efficient version of our attack also directly follows from the efficient version of Theorem 3.11.  $\square$

A natural Boolean cost function can be defined as

$$\text{cost}(h, \mathcal{X}) = \begin{cases} 0 & \text{if } h(d) = t \text{ for all } (d, t) \in \mathcal{X} \\ 1 & \text{otherwise} \end{cases}$$

where the cost function outputs 0 if the hypothesis is correct on all the examples in the test set. A special interesting case is where  $X'$  contains a single element  $t \leftarrow X$  sampled from  $X$  itself, but the adversary knows this test example and hopes to increase the error of classifying  $t$ .

**Corollary 4.8** (Doubling the error). *For small error  $\delta$ , we can double it by using  $p \approx 1/2$ . Indeed, for every deterministic learning algorithm  $L$  that outputs a hypothesis  $h$  by taking a sequence of  $n$  labeled examples generated by an oracle  $O_X$  and for every Boolean cost function  $\text{cost} : \mathcal{H} \times 2^{\mathcal{P}^c} \rightarrow \{0, 1\}$ , there exist a  $p$ -tampering training oracle  $\widehat{O}_X^p$ , using  $p = \frac{1}{2(1-\delta)}$ , that doubles the average cost  $\delta = \overline{\text{cost}}_{L^{O_X}}(\mathcal{X})$  into  $2\delta$ .*

## 5 Efficient $p$ -Tampering Attacks Biasing Bounded Functions

In this section we will formally prove Theorems 3.9 and 3.9. As described in Section 1.2, some of the ideas (and even notation) that we use here goes back to the original work of Austrin et. al [1] and here we show how to extend these arguments to the blockwise setting and overcome challenges that emerge.

Before doing so, we need to define some useful notation for the notions that naturally come up in our proofs. We will also make some basic observations about these quantities before proving our main theorems.

**Definition 5.1** (Functions  $\hat{f}, g, \mathcal{G}, \mathcal{A}, \mathcal{Q}$ ). Suppose  $f: \text{Supp}(\bar{X}) \mapsto \mathbb{R}$  is defined over a joint distribution  $\bar{X} = (X_1, \dots, X_n)$ ,  $i \in [n]$ , and  $x_{\leq i} \in \text{ValPref}(\bar{X})$  is a valid prefix for  $\bar{X}$ . Then we define the following with respect to  $f, \bar{X}, x_{\leq i}$ .

- $f_{x_{\leq i}}(\cdot)$  is a function defined as  $f_{x_{\leq i}}(x_{\geq i+1}) = f(x)$  where  $x = (x_{\leq i}, x_{\geq i+1})$ .
- $\hat{f}[x_{\leq i}] = \mathbb{E}_{x_{\geq i+1} \leftarrow (X_{\geq i+1} | x_{\leq i})} [f_{x_{\leq i}}(x_{\geq i+1})]$ . We also use  $\mu = \hat{f}[\emptyset]$  to denote  $\hat{f}[x_{\leq 0}] = \mathbb{E}[f(\bar{X})]$ .
- We define the *gain* of the “node”  $x_{\leq i}$  (compared to its parent  $x_{\leq i-1}$ ) as  $g[x_{\leq i}] = \hat{f}[x_{\leq i}] - \hat{f}[x_{\leq i-1}]$ . This defines the change in  $\hat{f}[x_{\leq i}]$  after moving to the  $i$ 'th block.
- For every  $x_{\leq i-1}$  and every distribution  $Z$  that *could depend* on  $x_{\leq i-1}$  (e.g.,  $Z$  is the output of a randomized algorithm that takes  $x_{\leq i-1}$  as input) and  $\text{Supp}(Z | x_{\leq i-1}) \subseteq \text{Supp}(X_i | x_{\leq i-1})$  we define:

- The *average of the gain* over the “children” of node  $x_{\leq i-1}$  under distribution  $(Z | x_{\leq i-1})$ :

$$\mathcal{G}_Z[x_{\leq i-1}] = \mathbb{E}_{x_i \leftarrow (Z | x_{\leq i-1})} [g[x_{\leq i}]].$$

- The average of the *absolute value* of the gains:

$$\mathcal{A}_Z[x_{\leq i-1}] = \mathbb{E}_{x_i \leftarrow (Z | x_{\leq i-1})} \left[ |g[x_{\leq i}]| \right].$$

- The average of the *squares* of the gains:

$$\mathcal{Q}_Z[x_{\leq i-1}] = \mathbb{E}_{x_i \leftarrow (Z | x_{\leq i-1})} \left[ g[x_{\leq i}]^2 \right].$$

**Notation.** Throughout following sections, whenever we define  $\bar{X}$  and  $f$ , then we will use all the notations defined in Definition 5.1 with respect to  $f$  and  $\bar{X}$  even if there are other distributions like  $\bar{Y}$  defined.

The following lemma directly follows from the definition of  $\mu$  and  $g[x_{\leq i}]$ .

**Proposition 5.2.** For every  $x \in \text{Supp}(\bar{X})$ ,  $f(x) = \mu + \sum_{i \in [n]} g[x_{\leq i}]$ .

The following two intuitive propositions also follow from the definition of  $\mathcal{G}_{X_i}[x_{\leq i-1}]$ .

**Proposition 5.3.** For every valid prefix  $x_{\leq i-1} \in \text{ValPref}(\bar{X})$ , we have  $\mathcal{G}_{X_i}[x_{\leq i-1}] = 0$ .

*Proof.*

$$\begin{aligned}
\mathcal{G}_{X_i}[x_{\leq i-1}] &= \mathbb{E}_{x_i \leftarrow (X_i | x_{\leq i-1})} g[x_{\leq i}] \\
(\text{by definition of } g[x_{\leq i}]) &= \mathbb{E}_{x_i \leftarrow (X_i | x_{\leq i-1})} [\hat{f}[x_{\leq i}] - \hat{f}[x_{\leq i-1}]] \\
(\hat{f}[x_{\leq i-1}] \text{ is independent of } x_i) &= \mathbb{E}_{x_i \leftarrow (X_i | x_{\leq i-1})} [\hat{f}[x_{\leq i}]] - \hat{f}[x_{\leq i-1}] \\
(\text{by definitions of } \hat{f}[x_{\leq i-1}], \hat{f}[x_{\leq i}]) &= \mathbb{E}_{x_i \leftarrow (X_i | x_{\leq i-1})} \left[ \mathbb{E}_{x_{\geq i+1} \leftarrow (X_{\geq i+1} | x_{\leq i})} [f(x)] \right] - \mathbb{E}_{x_{\geq i} \leftarrow (X_{\geq i} | x_{\leq i-1})} [f(x)] \\
&= \mathbb{E}_{x_{\geq i} \leftarrow (X_{\geq i} | x_{\leq i-1})} [f(x)] - \mathbb{E}_{x_{\geq i} \leftarrow (X_{\geq i} | x_{\leq i-1})} [f(x)] \\
&= \hat{f}[x_{\leq i-1}] - \hat{f}[x_{\leq i-1}] = 0.
\end{aligned}$$

□

**Proposition 5.4.** *Let  $f: \text{Supp}(\bar{X}) \mapsto \mathbb{R}$  be any real-output function. Then for any distribution  $\bar{Y}$  such that  $\text{Supp}(\bar{Y}) \subseteq \text{Supp}(\bar{X})$  it holds that  $\mathbb{E}[f(\bar{Y})] - \mathbb{E}[f(\bar{X})] = \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} [\mathcal{G}_{Y_i}[Y_{\leq i-1}]]$ .*

*Proof.*

$$\begin{aligned}
(\text{by Proposition 5.2}) \quad \mathbb{E}[f(\bar{Y})] - \mathbb{E}[f(\bar{X})] &= \mathbb{E}_{\bar{Y}} \left[ \sum_{i \in [n]} [g[x_{\leq i}]] \right] \\
(\text{by the linearity of expectation}) &= \sum_{i \in [n]} \mathbb{E}_{\bar{Y}} [g[Y_{\leq i}]] \\
&= \sum_{i \in [n]} \mathbb{E}_{(Y_{\leq i-1})} \mathbb{E}_{(Y_i | X_{\leq i-1} = Y_{\leq i-1})} [g[Y_{\leq i}]] \\
(\text{by the definition of } \mathcal{G}) &= \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} [\mathcal{G}_{Y_i}[Y_{\leq i-1}]].
\end{aligned}$$

□

The above analysis holds for any distribution  $\bar{Y}$  as long as  $\text{Supp}(\bar{Y}) \subseteq \text{Supp}(\bar{X})$ , but the following claim is about  $\rho$ -tampering variations.

**Proposition 5.5.** *For any probability tree  $\rho$  over  $\bar{X}$ , and any  $\rho$ -tampering variation  $\bar{Y}$  of  $\bar{X}$  generated by a (possibly randomized) tampering algorithm  $\text{Tam}$ , and for any  $y_{\leq i-1} \in \text{ValPref}(\bar{X})$ , it holds that  $\mathcal{G}_{Y_i}[y_{\leq i-1}] = \rho[y_{\leq i-1}] \cdot \mathcal{G}_{\text{Tam}}[y_{\leq i-1}]$ .*

*Proof.* The proof simply follows from the definition of  $\rho$ -tampering variations. When we sample from the distribution  $(Y_i | \bar{Y}_{\leq i-1} = y_{\leq i-1})$ , by definition, with probability  $1 - \rho[y_{\leq i-1}]$  we will be sampling  $Y_i$  from  $(X_i | X_{\leq i-1} = y_{\leq i-1})$  which by Proposition 5.3 leads to gaining  $\mathcal{G}_{X_i}[y_{\leq i-1}] = 0$ , and with probability  $\rho[y_{\leq i-1}]$  we will be sampling  $Y_i$  from  $\text{Tam}(y_{\leq i-1})$  which leads to gaining  $\mathcal{G}_{\text{Tam}}[y_{\leq i-1}]$ . Putting together, this implies an average gain of  $\rho[y_{\leq i-1}] \cdot \mathcal{G}_{\text{Tam}}[y_{\leq i-1}]$ . □

## 5.1 Biasing Real-Output Functions: Proving Theorem 3.9

In this Section we will prove our Theorem 3.9.

**Construction 5.6.** Let  $\bar{X} = (X_1, \dots, X_n)$  be the joint distribution and  $f: \text{Supp}(\bar{X}) \mapsto [-1, +1]$ . The *one rejection sampling* tampering algorithm ORSam works as follows. Given the valid prefix  $y_{\leq i-1} \in \text{ValPref}(\bar{X})$ , the tampering algorithm would sample  $y_{\geq i} \leftarrow (X_{\geq i} \mid y_{\leq i-1})$  by multiple invocations of the online sampler  $S$ . Then it computes  $s = f(y_1, \dots, y_n)$  and output from the following random variable.

$$T = \begin{cases} \text{Case 1: with probability } \frac{1+s}{2} \text{ output } y_i. \\ \text{Case 2: with probability } \frac{1-s}{2} \text{ output a fresh sample } y'_i \leftarrow S(y_{\leq i-1}). \end{cases}$$

**Claim 3.** Let  $f: \text{Supp}(\bar{X}) \rightarrow [-1, +1]$ ,  $\rho$  be every  $[p, q]$ -probability tree over  $\bar{X}$ , and  $\mu = \mathbb{E}[f(\bar{X})]$ . Then, the tampering algorithm ORSam of construction 5.6 generates a  $\rho$ -tampering variation  $\bar{Y}$  of  $\bar{X}$  such that

$$\mathbb{E}[f(\bar{Y})] \geq \mu + \frac{p}{3 + p \cdot \left( \frac{q}{1-q} \cdot (1 + |\mu|)^2 + 3 - 3\mu \right)} \cdot \text{Var}[f(\bar{X})].$$

For the special case that  $f: \text{Supp}(\bar{X}) \rightarrow \{+1, -1\}$  is Boolean, then we get a better bound of:

$$\mathbb{E}[f(\bar{Y})] \geq \mu + \frac{p}{2 + 2p} \cdot \text{Var}[f(\bar{X})]$$

We first prove Theorem 3.9 using Claim 3, and then we will prove Claim 3.

*Proof of Theorem 3.9.* We need to show that there is an attack that can bias  $f$  by  $\Omega(p) \cdot \text{Var}[f(\bar{X})]$ . For the Boolean case the proof follows directly from the statement of Claim 3. For the case of real-output functions we use an attacker that with probability 0.5 uses a fresh sample, and with probability 0.5 it runs the one-rejection sampling attack of Construction 5.6. This algorithm gives a  $\rho$ -tampering variation  $\bar{Y}$  of  $\bar{X}$  such that  $\forall y_{\leq i} \in \text{ValPref}(\bar{X})$ ,  $\frac{p}{2} \leq \rho[y_{\leq i}] \leq \frac{1}{2}$  so using Claim 3 we have:

$$\mathbb{E}[f(\bar{Y})] - \mathbb{E}[f(\bar{X})] \geq \frac{p/2}{3 + \frac{p}{2} \cdot ((1 + |\mu|)^2 - 3\mu + 3)} \cdot \text{Var}[f(\bar{X})] \geq \frac{p}{6 + 10p} \cdot \text{Var}[f(\bar{X})] \geq \Omega(p) \cdot \text{Var}[f(\bar{X})].$$

□

**Remark 5.7** (Bounds for special cases of real-valued functions). As shown above, Claim 3 can be used to derive the lower bound of  $\frac{p}{6+10p} \cdot \text{Var}[f(\bar{X})]$  for bias of general real-valued functions through a  $p$ -tampering. Now, for some natural special cases, we list the better bounds that are implied by Claim 3.

- **Balanced functions.** If the function  $f$  is balanced, or even  $\mu \geq 0$ , Claim 3 implies the stronger bound of  $\frac{p}{6+4p} \cdot \text{Var}[f(\bar{X})]$  for the bias.
- **Bounded tampering probability  $p$ .** If we know that the tampering probability does not happen with probability more than  $1/2$ , namely  $q \leq 1/2$  (e.g., if the tampering probability is fixed for some  $p = q \leq 1/2$ ), then we do not need to scale down the tampering probability and the one-resetting attack of Construction 5.6 gives bias at least  $\frac{p}{3+10p} \cdot \text{Var}[f(\bar{X})]$  already.
- **Having both properties.** If we have both of the properties above, we get the best of both which is a lower bound of  $\frac{p}{3+4p} \cdot \text{Var}[f(\bar{X})] \geq \frac{p}{7} \cdot \text{Var}[f(\bar{X})]$  for the achieved bias.

In the rest of this section, we will prove Claim 3. All along we use  $\bar{Y}$  to denote the  $\rho$ -tampering variation of  $\bar{X}$  generated by one rejection sampling algorithm ORSam of Construction 5.6.

**Claim 4.** *Let  $T \equiv \text{ORSam}(y_{\leq i-1})$  be a random variable defined over the randomness of ORSam running on a valid prefix  $y_{\leq i-1} \in \text{ValPref}(\bar{X})$ . The probability distribution of this random variable is:*

$$\Pr[T = y_i] = \left(1 + \frac{g[y_{\leq i}]}{2}\right) \cdot \Pr[X_i = y_i \mid y_{\leq i-1}].$$

*Proof.* We have two cases in the attack. We first compute the probability of Case 1.

$$\Pr[\text{Case 1} \wedge T = y_i] = \mathbb{E}_{y_{>i} \leftarrow (X_{>i} | y_{\leq i-1})} \left[ \frac{1 + f(y)}{2} \right] \cdot \Pr[X_i = y_i \mid y_{\leq i-1}] = \left( \frac{1 + \hat{f}[y_{\leq i}]}{2} \right) \cdot \Pr[X_i = y_i \mid y_{\leq i-1}].$$

On the other hand, the probability of Case 2 is

$$\begin{aligned} \Pr[\text{Case 2} \wedge T = y_i] &= \Pr[T = y_i \mid \text{Case 2}] \cdot \Pr[\text{Case 2}] \\ &= \Pr[X_i = y_i \mid y_{\leq i-1}] \cdot \mathbb{E}_{y_{>i-1} \leftarrow (X_{>i-1} | y_{\leq i-1})} \left[ \frac{1 - f(y)}{2} \right] \\ &= \Pr[X_i = y_i \mid y_{\leq i-1}] \cdot \left( \frac{1 - \hat{f}[y_{\leq i-1}]}{2} \right). \end{aligned}$$

Thus, we have

$$\begin{aligned} \Pr[T = y_i] &= \Pr[\text{Case 1} \wedge T = y_i] + \Pr[\text{Case 2} \wedge T = y_i] \\ &= \left( \frac{1 + \hat{f}[y_{\leq i}]}{2} \right) \cdot \Pr[X_i = y_i \mid y_{\leq i-1}] + \Pr[X_i = y_i \mid y_{\leq i-1}] \cdot \left( \frac{1 - \hat{f}[y_{\leq i-1}]}{2} \right) \\ &= \left( 1 + \frac{g[y_{\leq i}]}{2} \right) \cdot \Pr[X_i = y_i \mid y_{\leq i-1}]. \end{aligned}$$

□

**Corollary 5.8.** *For any  $y_{\leq i} \in \text{ValPref}(\bar{X})$ , it holds that*

$$\Pr[Y_i = y_i \mid y_{\leq i-1}] = \left( 1 + \frac{\rho[y_{\leq i-1}] \cdot g[y_{\leq i}]}{2} \right) \cdot \Pr[X_i = y_i \mid y_{\leq i-1}].$$

*Proof.* By definition of  $\bar{Y}$  we have

$$\begin{aligned} \Pr[Y_i = y_i \mid y_{\leq i-1}] &= (1 - \rho[y_{\leq i-1}]) \cdot \Pr[X_i = y_i \mid y_{\leq i-1}] + \rho[y_{\leq i-1}] \cdot \Pr[y_i = \text{ORSam}(y_{\leq i-1})] \\ \text{(by Claim 4)} \quad &= (1 - \rho[y_{\leq i-1}] + \rho[y_{\leq i-1}] \cdot (1 + \frac{g[y_{\leq i}]}{2})) \Pr[X_i = y_i \mid y_{\leq i-1}] \\ &= \left( 1 + \frac{\rho[y_{\leq i-1}] \cdot g[y_{\leq i}]}{2} \right) \cdot \Pr[X_i = y_i \mid y_{\leq i-1}]. \end{aligned}$$

□

**Lemma 1.** *Let  $\bar{X} = (X_1, \dots, X_n)$ . For every function  $f: \text{Supp}(\bar{X}) \rightarrow [-1, +1]$  and every  $[p, 1]$ -probability tree  $\rho$  over  $\bar{X}$ , if  $\bar{Y}$  is the  $\rho$ -tampering variation of distribution  $\bar{X}$  generated by tampering algorithm ORSam of construction 5.6, and if  $\mu = \mathbb{E}[f(\bar{X})]$ , then it holds that*

$$\mathbb{E}[f(\bar{Y})] \geq \mu + \frac{p}{2(1+p)} \cdot (\mathbb{E}[f(\bar{Y})^2] - \mu^2).$$



**Boolean case.** The above finishes the proof for the case of Boolean  $f$ , because  $f(\bar{Y})^2 = 1$ , and so

$$\mathbb{E}[f(\bar{Y})] - \mu \geq \frac{p}{2(1+p)} \cdot (\mathbb{E}[f(\bar{Y})^2] - \mu^2) = \frac{p}{2(1+p)} \cdot (1 - \mu^2) = \frac{p}{2(1+p)} \cdot \text{Var}[f(\bar{X})].$$

Before proving the above lemma, we will need to prove several other claims.

**Claim 5** (One rejection sampling's local gains). *For any  $y_{\leq i} \in \text{ValPref}(\bar{X})$ , it holds that*

$$\mathcal{G}_{\text{ORSam}}[y_{\leq i-1}] = \mathcal{Q}_{X_i}[y_{\leq i-1}]/2.$$

*Proof.* First note that  $\mathcal{G}_{\text{ORSam}}[y_{\leq i-1}] = \sum_{y_i} \Pr[y_i = \text{ORSam}(y_{\leq i})] \cdot g[y_{\leq i}]$ . By Claim 4 we get

$$\begin{aligned} \mathcal{G}_{\text{ORSam}}[y_{\leq i-1}] &= \sum_{y_i} \Pr[X_i = y_i \mid y_{\leq i-1}] \cdot \left(1 + \frac{g[y_{\leq i}]}{2}\right) \cdot g[y_{\leq i}] \\ &= \sum_{y_i} \Pr[X_i = y_i \mid y_{\leq i-1}] \cdot g[y_{\leq i}] + \sum_{y_i} \Pr[X_i = y_i \mid y_{\leq i-1}] \cdot \frac{g[y_{\leq i}]^2}{2} \\ &= \mathcal{G}_{X_i}[y_{\leq i-1}] + \frac{\mathcal{Q}_{X_i}[y_{\leq i-1}]}{2}. \end{aligned}$$

By Proposition 5.3 we also know that  $\mathcal{G}_{X_i}[y_{\leq i-1}] = 0$ , so  $\mathcal{G}_{\text{ORSam}}[y_{\leq i-1}] = \mathcal{Q}_{X_i}[y_{\leq i-1}]/2$ .  $\square$

**Corollary 5.9.** *For any  $y_{\leq i-1} \in \text{ValPref}(\bar{X})$ , it holds that  $\mathcal{G}_{Y_i}[y_{\leq i-1}] = \frac{\rho[y_{\leq i-1}]}{2} \cdot \mathcal{Q}_{X_i}[y_{\leq i-1}]$ .*

*Proof.*

$$\begin{aligned} \mathcal{G}_{Y_i}[y_{\leq i-1}] &= \sum_{y_i} \Pr[y_i = Y_i \mid y_{\leq i-1}] \cdot g[y_{\leq i}] \\ &= \sum_{y_i} \left( (1 - \rho[y_{\leq i-1}]) \cdot \Pr[y_i = X_i \mid y_{\leq i-1}] \right) \cdot g[y_{\leq i}] \\ &\quad + \sum_{y_i} \left( \rho[y_{\leq i-1}] \cdot \Pr[y_i = \text{ORSam}(y_{\leq i-1})] \right) \cdot g[y_{\leq i}] \\ &= (1 - \rho[y_{\leq i-1}]) \cdot \mathcal{G}_{X_i}[y_{\leq i-1}] + \rho[y_{\leq i-1}] \cdot \mathcal{G}_{\text{ORSam}}[y_{\leq i-1}] \\ \text{(by Proposition 5.3)} &= \rho[y_{\leq i-1}] \cdot \mathcal{G}_{\text{ORSam}}[y_{\leq i-1}] \\ \text{(by Claim 5)} &= \frac{\rho[y_{\leq i-1}]}{2} \cdot \mathcal{Q}_{X_i}[y_{\leq i-1}]. \end{aligned}$$

$\square$

**Corollary 5.10.**  $\mathbb{E}_{\bar{Y}}[f(\bar{Y})] = \mu + \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} \left[ \frac{\rho[Y_{\leq i-1}]}{2} \cdot \mathcal{Q}_{X_i}[Y_{\leq i-1}] \right]$ .

*Proof.* Using Claim 5.4, we have  $\mathbb{E}_{\bar{Y}}[f(\bar{Y})] = \mu + \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} [\mathcal{G}_{Y_i}[Y_{\leq i-1}]]$ . By also using Corollary 5 we obtain  $\mathbb{E}_{\bar{Y}}[f(\bar{Y})] = \mu + \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} \left[ \frac{\rho[Y_{\leq i-1}]}{2} \cdot \mathcal{Q}_{X_i}[Y_{\leq i-1}] \right]$ .  $\square$

**Claim 6.** *For every  $x \in \text{Supp}(\bar{X})$ , it holds that*

$$f(x)^2 = \mu^2 + \sum_{i=1}^n \left( g[x_{\leq i}]^2 + 2\hat{f}[x_{\leq i-1}] \cdot g[x_{\leq i}] \right).$$

*Proof.* By squaring the equation in Proposition 5.2 we get

$$f(x)^2 = \mu^2 + \sum_{i=1}^n g[x_{\leq i}]^2 + 2 \sum_{i=1}^n g[x_{\leq i}] \cdot \left( \mu + \sum_{j=1}^{i-1} g[x_{\leq j}] \right)$$

By the definition of  $g[x_{\leq j}]$  it holds that  $\hat{f}[x_{\leq i-1}] = \mu + \sum_{j=1}^{i-1} g[x_{\leq j}]$ . So we get

$$f(x)^2 = \mu^2 + \sum_{i=1}^n \left( g[x_{\leq i}]^2 + 2\hat{f}[x_{\leq i-1}] \cdot g[x_{\leq i}] \right).$$

□

**Claim 7.** For any  $y_{\leq i-1} \in \text{ValPref}(\bar{X})$ , it holds that

$$\mathcal{Q}_{Y_i}[y_{\leq i-1}] = \mathcal{Q}_{X_i}[y_{\leq i-1}] + \mathbb{E}_{X_i|y_{\leq i-1}} \left[ \frac{\rho[y_{\leq i-1}]}{2} \cdot g[(y_{\leq i-1}, X_i)]^3 \right].$$

*Proof.*

$$\begin{aligned} \mathcal{Q}_{Y_i}[y_{\leq i-1}] &= \mathbb{E}_{Y_i|y_{\leq i-1}} [g[(y_{\leq i-1}, Y_i)]^2] \\ &= \sum_{y_i} \Pr[Y_i = y_i | y_{\leq i-1}] \cdot g[y_{\leq i}]^2 \\ \text{(by Corollary 5.8)} &= \sum_{y_i} \left( 1 + \frac{\rho[y_{\leq i-1}]}{2} \cdot g[y_{\leq i}] \right) \cdot \Pr[X_i = y_i | y_{\leq i-1}] \cdot g[y_{\leq i}]^2 \\ &= \sum_{y_i} \Pr[X_i = y_i | y_{\leq i-1}] \cdot g[y_{\leq i}]^2 + \frac{\rho[y_{\leq i-1}]}{2} \cdot \Pr[X_i = y_i | y_{\leq i-1}] \cdot g[y_{\leq i}]^3 \\ &= \mathcal{Q}_{X_i}[y_{\leq i-1}] + \mathbb{E}_{X_i|y_{\leq i-1}} \left[ \frac{\rho[y_{\leq i-1}]}{2} g[(y_{\leq i-1}, X_i)]^3 \right] \end{aligned}$$

□

**Claim 8.** For any  $y_{\leq i-1} \in \text{ValPref}(\bar{X})$ , it holds that

$$\hat{f}[y_{\leq i-1}] \cdot \mathcal{Q}_{X_i}[y_{\leq i-1}] + \mathbb{E}_{X_i|y_{\leq i-1}} [g[(y_{\leq i-1}, X_i)]^3] \leq \mathcal{Q}_{X_i}[y_{\leq i-1}].$$

*Proof.*  $\mathbb{E}_{X_i|y_{\leq i-1}} [g[(y_{\leq i-1}, X_i)]^3] + \hat{f}[y_{\leq i-1}] \cdot \mathcal{Q}_{X_i}[y_{\leq i-1}]$  is equal to:

$$\begin{aligned} &= \mathbb{E}_{X_i|y_{\leq i-1}} [g[(y_{\leq i-1}, X_i)]^3] + \hat{f}[y_{\leq i-1}] \cdot \mathbb{E}_{X_i|y_{\leq i-1}} [g[(y_{\leq i-1}, X_i)]^2] \\ &= \mathbb{E}_{X_i|y_{\leq i-1}} [g[(y_{\leq i-1}, X_i)]^3] + \mathbb{E}_{X_i|y_{\leq i-1}} [\hat{f}[y_{\leq i-1}] \cdot g[(y_{\leq i-1}, X_i)]^2] \\ \text{(by linearity of expectations)} &= \mathbb{E}_{X_i|y_{\leq i-1}} \left[ g[(y_{\leq i-1}, X_i)]^3 + \hat{f}[y_{\leq i-1}] \cdot g[(y_{\leq i-1}, X_i)]^2 \right] \\ &= \mathbb{E}_{X_i|y_{\leq i-1}} \left[ \left( g[(y_{\leq i-1}, X_i)] + \hat{f}[y_{\leq i-1}] \right) \cdot g[(y_{\leq i-1}, X_i)]^2 \right] \\ \text{(by the definition of } g[(y_{\leq i-1}, X_i)] \text{)} &= \mathbb{E}_{X_i|y_{\leq i-1}} [\hat{f}[(y_{\leq i-1}, X_i)] \cdot g[(y_{\leq i-1}, X_i)]^2]. \end{aligned}$$

Now, because  $\hat{f}[(y_{\leq i-1}, X_i)] \leq 1$ , the above is at most  $\mathbb{E}_{X_i|y_{\leq i-1}} [g[(y_{\leq i-1}, X_i)]^2] = \mathcal{Q}_{X_i}[y_{\leq i-1}]$ . □

**Claim 9.** For any  $[p, 1]$ -probability tree  $\rho$  over  $\bar{X}$  it holds that

$$\mathbb{E}[f(\bar{Y})^2] \leq \mu^2 + \frac{1+p}{p} \cdot \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} [\rho[Y_{\leq i-1}] \mathcal{Q}_{X_i}[Y_{\leq i-1}]].$$

*Proof.* Using Claim 6 we have

$$\begin{aligned} \frac{\mathbb{E}[f(\bar{Y})^2]}{\bar{Y}} - \mu^2 &= \sum_{i=1}^n \mathbb{E}_{\bar{Y}} \left[ g[Y_{\leq i}]^2 + 2\hat{f}[Y_{\leq i-1}] \cdot g[Y_{\leq i}] \right] \\ &= \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} \left[ \mathbb{E}_{Y_i|Y_{\leq i-1}} [g[Y_{\leq i}]^2] \right] + 2 \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} [\hat{f}[Y_{\leq i-1}] \cdot \mathbb{E}_{Y_i|Y_{\leq i-1}} [g[Y_{\leq i}]]] \\ &= \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} [\mathcal{Q}_{Y_i}[Y_{\leq i-1}]] + 2 \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} [\hat{f}[Y_{\leq i-1}] \cdot \mathcal{G}_{Y_i}[\bar{Y}_{\leq i-1}]] \\ \text{(by Claim 7)} &= \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} \left[ \mathcal{Q}_{X_i}[Y_{\leq i-1}] + \frac{\rho[Y_{\leq i-1}]}{2} \mathbb{E}_{X_i|Y_{\leq i-1}} [g[(Y_{\leq i-1}, X_i)]^3] \right] \\ &\quad + \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} \left[ 2\hat{f}[Y_{\leq i-1}] \cdot \mathcal{G}_{Y_i}[Y_{\leq i-1}] \right] \\ \text{(by Corollary 5.9)} &= \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} \left[ \mathcal{Q}_{X_i}[Y_{\leq i-1}] + \frac{\rho[Y_{\leq i-1}]}{2} \mathbb{E}_{X_i|Y_{\leq i-1}} [g[(Y_{\leq i-1}, X_i)]^3] \right] \\ &\quad + \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} \left[ \rho[Y_{\leq i-1}] \hat{f}[Y_{\leq i-1}] \mathcal{Q}_{X_i}[Y_{\leq i-1}] \right] \\ \text{(by Claim 8)} &\leq \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} \left[ \mathcal{Q}_{X_i}[Y_{\leq i-1}] + \frac{\rho[Y_{\leq i-1}]}{2} \cdot \mathcal{Q}_{X_i}[Y_{\leq i-1}] \right] \\ &\quad + \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} \left[ \frac{\rho[Y_{\leq i-1}]}{2} \cdot \hat{f}[Y_{\leq i-1}] \cdot \mathcal{Q}_{X_i}[Y_{\leq i-1}] \right] \\ \text{(by } \hat{f}[Y_{\leq i-1}] \leq 1) &\leq \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} [(1 + \rho[Y_{\leq i-1}]) \cdot \mathcal{Q}_{X_i}[Y_{\leq i-1}]] \\ \text{(by } \rho[Y_{\leq i-1}] \geq p) &\leq \left( \frac{1}{p} + 1 \right) \cdot \sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} [\rho[Y_{\leq i-1}] \cdot \mathcal{Q}_{X_i}[Y_{\leq i-1}]]. \end{aligned}$$

□

Now we will prove Lemma 1.

*Proof of Lemma 1.* Using Claim 9 we have

$$\sum_{i=1}^n \mathbb{E}_{Y_{\leq i-1}} [\rho[Y_{\leq i-1}] \cdot \mathcal{Q}_{X_i}[Y_{\leq i-1}]] \geq \frac{p}{1+p} \cdot \frac{\mathbb{E}[f(\bar{Y})^2]}{\bar{Y}} - \mu^2$$

By also applying Corollary 5.10 we get  $\mathbb{E}[f(\bar{Y})] \geq \mu + \frac{p}{2(1+p)} \cdot (\mathbb{E}[f(\bar{Y})^2] - \mu^2)$ .

□

**Lemma 2.** Let  $f: \bar{X} \rightarrow [-1, +1]$ ,  $\mu = \mathbb{E}[f(\bar{X})]$ ,  $\nu = \text{Var}[f(\bar{X})]$ , and  $\rho$  be a  $[0, q]$ -probability tree over  $\bar{X}$ . If  $\bar{Y}$  is the  $\rho$ -tampering variation of distribution  $\bar{X}$  generated by ORSam of construction 5.6, then:

$$\mathbb{E}_{\bar{Y}}[f(\bar{Y})] + \frac{3(1-q)}{2q \cdot (1+|\mu|)^2} \cdot \mathbb{E}_{\bar{Y}} [(f(\bar{Y}) - \mu)^2] \geq \mu + \frac{(1-q) \cdot \nu}{q \cdot (1+|\mu|)^2}.$$

Before proving Lemma 2 we need to define a few useful functions.

**Definition 5.11** (Potential function). Let  $t: \text{Supp}(\bar{X}) \rightarrow [-1, +1]$  be an arbitrary function also let  $\hat{t}$  be defined as  $\hat{f}$  of Definition 5.1, namely,  $\hat{t}[x_{\leq i}] = \mathbb{E}_{x_{\geq i+1} \leftarrow (X_{\geq i+1}|x_{\leq i})}[t_{x_{\leq i}}(x_{\geq i+1})]$ . We define the potential function  $\Phi: \text{ValPref}(\bar{X}) \rightarrow \mathbb{R}$  based on  $t$  as follows

$$\Phi(y_{\leq i}) = \hat{f}[y_{\leq i}] + \frac{1-q}{q} \cdot \hat{t}[y_{\leq i}] + \frac{1-q}{2q} \cdot (\hat{t}[y_{\leq i}])^2.$$

**Claim 10** (Potential function does not decrease).  $\mathbb{E}[\Phi(Y_{\leq i})] \geq \mathbb{E}[\Phi(Y_{\leq i-1})]$ .

Before proving Claim 10, note that Lemma 2 immediately follows from Claim 10.

*Proof of Lemma 2.* Using Claim 10 together with a simple induction we get

$$\mathbb{E}[\Phi(Y_{\leq n})] \geq \mathbb{E}[\Phi(Y_{\leq 0})].$$

Now by setting  $t(y) = \left(\frac{f(y)-\mu}{1+|\mu|}\right)^2$ , which its outputs are in range  $[0, 1]$ , on the one side, we get

$$\begin{aligned} \mathbb{E}[\Phi(Y_{\leq n})] &= \mathbb{E}_{\bar{Y}}[f(\bar{Y})] + \frac{1-q}{q} \cdot \mathbb{E}_{\bar{Y}} \left[ \left( \frac{f(\bar{Y}) - \mu}{1+|\mu|} \right)^2 \right] + \frac{1-q}{2q} \cdot \mathbb{E}_{\bar{Y}} \left[ \left( \frac{f(\bar{Y}) - \mu}{1+|\mu|} \right)^4 \right] \\ &\leq \mathbb{E}_{\bar{Y}}[f(\bar{Y})] + \frac{3(1-q)}{2q} \cdot \mathbb{E}_{\bar{Y}} \left[ \left( \frac{f(\bar{Y}) - \mu}{1+|\mu|} \right)^2 \right] \end{aligned}$$

and, on the other side, by letting  $i = 0$ , we get

$$\begin{aligned} \mathbb{E}[\Phi(Y_{\leq 0})] &= \mathbb{E}_{\bar{X}}[f(\bar{X})] + \frac{1-q}{q} \cdot \mathbb{E}_{\bar{X}} \left[ \left( \frac{f(\bar{X}) - \mu}{1+|\mu|} \right)^2 \right] + \frac{1-q}{2q} \cdot \mathbb{E}_{\bar{X}} \left[ \left( \frac{f(\bar{X}) - \mu}{1+|\mu|} \right)^4 \right] \\ &= \mu + \frac{(1-q) \cdot \nu}{q \cdot (1+|\mu|)^2} + \frac{(1-q) \cdot \nu^2}{2q \cdot (1+|\mu|)^4} \\ &\geq \mu + \frac{(1-q) \cdot \nu}{q \cdot (1+|\mu|)^2}. \end{aligned}$$

□

Now, we prove Claim 10.

*Proof of Claim 10.* For every  $y_{\leq i} \in \text{ValPref}(\bar{X})$  we define  $r[y_{\leq i}] = \hat{t}[y_{\leq i}] - \hat{t}[y_{\leq i-1}]$ . It is easy to see that  $\mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})}[r[y_{\leq i}]] = 0$ , based on Proposition 5.3. Now we have

$$\begin{aligned}
\mathbb{E}[\Phi(Y_{\leq i})] &= \mathbb{E}_{Y_{\leq i}} [\hat{f}[Y_{\leq i}]] + \frac{1-q}{q} \mathbb{E}_{Y_{\leq i}} [\hat{t}[Y_{\leq i}]] + \frac{1-q}{2q} \mathbb{E}_{Y_{\leq i}} [\hat{t}[Y_{\leq i}]^2] \\
&= \mathbb{E}_{Y_{\leq i-1}} \left[ \hat{f}[Y_{\leq i}] + \mathbb{E}_{Y_i | Y_{\leq i-1}} [g[Y_{\leq i}]] \right] + \frac{1-q}{q} \mathbb{E}_{Y_{\leq i-1}} \left[ \hat{t}[Y_{\leq i-1}] + \mathbb{E}_{Y_i | Y_{\leq i-1}} [r[Y_{\leq i}]] \right] \\
&\quad + \frac{1-q}{2q} \mathbb{E}_{Y_{\leq i-1}} \left[ \hat{t}[Y_{\leq i-1}]^2 + \mathbb{E}_{Y_i | Y_{\leq i-1}} [2 \cdot r[Y_{\leq i}] \cdot \hat{t}[Y_{\leq i-1}] + r[Y_{\leq i}]^2] \right] \\
&= \mathbb{E}_{Y_{\leq i-1}} [\hat{f}[Y_{\leq i}]] + \frac{1-q}{q} \mathbb{E}_{Y_{\leq i-1}} [\hat{t}[Y_{\leq i-1}]] + \frac{1-q}{2q} \mathbb{E}_{Y_{\leq i-1}} [\hat{t}[Y_{\leq i-1}]^2] \\
&\quad + \mathbb{E}_{Y_{\leq i}} [g[Y_{\leq i}]] + \frac{1-q}{q} \mathbb{E}_{Y_{\leq i}} [r[Y_{\leq i}]] + \frac{1-q}{2q} \mathbb{E}_{Y_{\leq i}} [2 \cdot r[Y_{\leq i}] \cdot \hat{t}[Y_{\leq i-1}] + r[Y_{\leq i}]^2] \\
&= \mathbb{E}[\Phi(Y_{\leq i-1})] + \mathbb{E}_{Y_{\leq i}} [g[Y_{\leq i}]] + \frac{1-q}{q} \mathbb{E}_{Y_{\leq i}} [r[Y_{\leq i}]] \\
&\quad + \frac{1-q}{2q} \mathbb{E}_{Y_{\leq i}} [2 \cdot r[Y_{\leq i}] \cdot \hat{t}[Y_{\leq i-1}] + r[Y_{\leq i}]^2] \\
&= \mathbb{E}[\Phi(Y_{\leq i-1})] + \mathbb{E}_{Y_{\leq i-1}} \left[ \mathbb{E}_{Y_i | Y_{\leq i-1}} [g[Y_{\leq i}]] + \frac{1-q}{q} \mathbb{E}_{Y_i | Y_{\leq i-1}} [r[Y_{\leq i}]] \right] \\
&\quad + \mathbb{E}_{Y_{\leq i-1}} \left[ \frac{1-q}{2q} \mathbb{E}_{Y_i | Y_{\leq i-1}} [2\hat{t}[Y_{\leq i-1}]r[Y_{\leq i}] + r[Y_{\leq i}]^2] \right].
\end{aligned}$$

Now, it is enough to show that the the following quantity  $c$  is positive for all valid  $y_{\leq i-1} \in \text{ValPref}(\bar{X})$ :

$$c[y_{\leq i-1}] = \mathbb{E}_{y_i \leftarrow (Y_i | y_{\leq i-1})} \left[ g[y_{\leq i}] + \frac{1-q}{q} \cdot r[y_{\leq i}] + \frac{1-q}{2q} \cdot (2\hat{t}[y_{\leq i-1}] \cdot r[y_{\leq i}] + r[y_{\leq i}]^2) \right].$$

We decompose  $c$  into four value  $c_1, c_2, c_3, c_4$  as follows:

$$\begin{aligned}
c_1[y_{\leq i-1}] &= \mathbb{E}_{y_i \leftarrow (Y_i | y_{\leq i-1})} [g[y_{\leq i}]], & c_2[y_{\leq i-1}] &= \mathbb{E}_{y_i \leftarrow (Y_i | y_{\leq i-1})} [r[y_{\leq i}]] \\
c_3[y_{\leq i-1}] &= \mathbb{E}_{y_i \leftarrow (Y_i | y_{\leq i-1})} [2\hat{t}[y_{\leq i-1}] \cdot r[y_{\leq i}]], & c_4[y_{\leq i-1}] &= \mathbb{E}_{y_i \leftarrow (Y_i | y_{\leq i-1})} [r[y_{\leq i}]^2].
\end{aligned}$$

Based on these definitions we have  $c = c_1 + \frac{1-q}{q} \cdot c_2 + \frac{1-q}{2q} \cdot c_3 + \frac{1-q}{2q} \cdot c_4$ . We bound each one of these values separately then add them together. For  $c_1$  we have,

$$\begin{aligned}
c_1[y_{\leq i-1}] &= \mathbb{E}_{y_i \leftarrow (Y_i | y_{\leq i-1})} [g[y_{\leq i}]] = \sum_{y_i} \Pr[y_i = Y_i | y_{\leq i-1}] \cdot g[y_{\leq i}] \\
\text{(by Corollary 5.8)} &= \sum_{y_i} \left( 1 + \frac{\rho[y_{\leq i-1}] \cdot g[y_{\leq i}]}{2} \right) \cdot \Pr[X_i = y_i | y_{\leq i-1}] \cdot g[y_{\leq i}] \\
&= \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} [g[y_{\leq i}]] + \frac{\rho[y_{\leq i-1}]}{2} \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} [g[y_{\leq i}]^2] \\
\text{(by Proposition 5.3)} &= \frac{\rho[y_{\leq i-1}]}{2} \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} [g[y_{\leq i}]^2].
\end{aligned}$$

For  $c_2$  we have,

$$\begin{aligned}
c_2[y_{\leq i-1}] &= \mathbb{E}_{y_i \leftarrow (Y_i | y_{\leq i-1})} [r[y_{\leq i}]] = \sum_{y_i} \Pr[y_i = Y_i | y_{\leq i-1}] \cdot r[y_{\leq i}] \\
(\text{by Corollary 5.8}) &= \sum_{y_i} \left( 1 + \frac{\rho[y_{\leq i-1}] \cdot g[y_{\leq i}]}{2} \right) \cdot \Pr[X_i = y_i | y_{\leq i-1}] \cdot r[y_{\leq i}] \\
&= \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} [r[y_{\leq i}]] + \frac{\rho[y_{\leq i-1}]}{2} \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} [g[y_{\leq i}] \cdot r[y_{\leq i}]] \\
(\text{by Proposition 5.3}) &= \frac{\rho[y_{\leq i-1}]}{2} \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} [g[y_{\leq i}] \cdot r[y_{\leq i}]] \\
&\geq \frac{\rho[y_{\leq i-1}]}{2} \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} [-|g[y_{\leq i}]| \cdot |r[y_{\leq i}]|].
\end{aligned}$$

For  $c_3$  we have,

$$\begin{aligned}
c_3[y_{\leq i-1}] &= \mathbb{E}_{y_i \leftarrow (y_i | y_{\leq i-1})} [2\hat{t}[y_{\leq i-1}] \cdot r[y_{\leq i}]] \\
&= \sum_{y_i} 2 \Pr[y_i = Y_i | y_{\leq i-1}] \cdot \hat{t}[y_{\leq i-1}] \cdot r[y_{\leq i}] \\
(\text{by Corollary 5.8}) &= \sum_{y_i} 2 \left( 1 + \frac{\rho[y_{\leq i-1}] \cdot g[y_{\leq i}]}{2} \right) \cdot \Pr[X_i = y_i | y_{\leq i-1}] \cdot 2\hat{t}[y_{\leq i-1}] \cdot r[y_{\leq i}] \\
&= \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} [2\hat{t}[y_{\leq i-1}] \cdot r[y_{\leq i}]] + \frac{\rho[y_{\leq i-1}]}{2} \cdot \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} [2\hat{t}[y_{\leq i-1}] \cdot r[y_{\leq i}] \cdot g[y_{\leq i}]] \\
&= 2\hat{t}[y_{\leq i-1}] \cdot \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} [r[y_{\leq i}]] + \rho[y_{\leq i-1}] \cdot \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} [\hat{t}[y_{\leq i-1}] \cdot r[y_{\leq i}] \cdot g[y_{\leq i}]] \\
(\text{by Proposition 5.3}) &= \rho[y_{\leq i-1}] \cdot \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} [\hat{t}[y_{\leq i-1}] \cdot r[y_{\leq i}] \cdot g[y_{\leq i}]] \\
(|\hat{t}[y_{\leq i-1}]| \text{ is at most } 1) &\geq \rho[y_{\leq i-1}] \cdot \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} [-|g[y_{\leq i}]| \cdot |r[y_{\leq i}]|].
\end{aligned}$$

And finally for  $c_4$  we have,

$$\begin{aligned}
c_4[y_{\leq i-1}] &= \mathbb{E}_{y_i \leftarrow (Y_i | y_{\leq i-1})} [r[y_{\leq i}]^2] = \sum_{y_i} \Pr[y_i = Y_i | y_{\leq i-1}] \cdot r[y_{\leq i}]^2 \\
(\text{by Corollary 5.8}) &= \sum_{y_i} \left( 1 + \frac{\rho[y_{\leq i-1}] \cdot g[y_{\leq i}]}{2} \right) \cdot \Pr[X_i = y_i | y_{\leq i-1}] \cdot r[y_{\leq i}]^2 \\
&= \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} \left[ \left( 1 + \frac{\rho[y_{\leq i-1}] \cdot g[y_{\leq i}]}{2} \right) \cdot r[y_{\leq i}]^2 \right] \\
(g[y_{\leq i}] \text{ is at least } -2) &\geq (1 - \rho[y_{\leq i-1}]) \cdot \mathbb{E}_{y_i \leftarrow (X_i | y_{\leq i-1})} [r[y_{\leq i}]^2].
\end{aligned}$$

By adding this inequalities together we have

$$\begin{aligned}
c[y_{\leq i-1}] &\geq \mathbb{E}_{y_i \leftarrow Y_i | y_{\leq i-1}} \left[ \frac{\rho[y_{\leq i-1}]}{2} g[y_{\leq i}]^2 - \frac{\rho[y_{\leq i-1}](1-q)}{q} |g[y_{\leq i}]| \cdot |r[y_{\leq i}]| \right] \\
&\quad + \mathbb{E}_{y_i \leftarrow Y_i | y_{\leq i-1}} \left[ \left( \frac{(1-q)(1-\rho[y_{\leq i-1}])}{2q} r[y_{\leq i}]^2 \right) \right] \\
&= \mathbb{E}_{y_i \leftarrow Y_i | y_{\leq i-1}} \left[ \frac{\rho[y_{\leq i-1}]}{2} \left( g[y_{\leq i}]^2 - \frac{2(1-q)}{q} |g[y_{\leq i}]| \cdot |r[y_{\leq i}]| \right) \right] \\
&\quad + \mathbb{E}_{y_i \leftarrow Y_i | y_{\leq i-1}} \left[ \frac{\rho[y_{\leq i-1}]}{2} \left( \frac{(1-q) \cdot (1-\rho[y_{\leq i-1}])}{q \cdot \rho[y_{\leq i-1}]} r[y_{\leq i}]^2 \right) \right] \\
(\text{by } \rho[y_{\leq i-1}] \leq q) &\geq \mathbb{E}_{y_i \leftarrow Y_i | y_{\leq i-1}} \left[ \frac{\rho[y_{\leq i-1}]}{2} \left( g[y_{\leq i}]^2 - \frac{2(1-q)}{q} |g[y_{\leq i}]| \cdot |r[y_{\leq i}]| \right) \right] \\
&\quad + \mathbb{E}_{y_i \leftarrow Y_i | y_{\leq i-1}} \left[ \frac{\rho[y_{\leq i-1}]}{2} \left( \frac{(1-q)^2}{q^2} r[y_{\leq i}]^2 \right) \right] \\
&= \mathbb{E}_{y_i \leftarrow Y_i | y_{\leq i-1}} \left[ \frac{\rho[y_{\leq i-1}]}{2} \left( |g[y_{\leq i}]| - \left| \frac{(1-q) \cdot r[y_{\leq i}]}{q} \right| \right)^2 \right] \geq 0.
\end{aligned}$$

So  $c[y_{\leq i-1}] \geq 0$  for every valid prefix  $y_{\leq i-1} \in \text{ValPref}(\bar{X})$  which finishes the proof of Claim 10.  $\square$

Finally, we prove Claim 3.

*Proof of Claim 3.* Let  $\mu' = \mathbb{E}[f(\bar{Y})]$  and  $b = \mu' - \mu$ , we have

$$\begin{aligned}
\mathbb{E}[f(\bar{Y})] + \frac{3(1-q)}{2q \cdot (1+|\mu|)^2} \cdot \mathbb{E}[(f(\bar{Y}) - \mu)^2] &= \mu' + \frac{3(1-q)}{2q \cdot (1+|\mu|)^2} \cdot (\mathbb{E}[f(\bar{Y})^2] + \mu^2 - 2\mu \cdot \mu') \\
(\text{By } \mu' = \mu + b) &= \mu' + \frac{3(1-q)}{2q \cdot (1+|\mu|)^2} \cdot (\mathbb{E}[f(\bar{Y})^2] - \mu^2 - 2\mu \cdot b) \\
(\text{By Lemma 1}) &\leq \mu' + \frac{3(1-q)}{2q \cdot (1+|\mu|)^2} \cdot \left( \frac{b \cdot (2+2p)}{p} - 2\mu \cdot b \right)
\end{aligned}$$

Now using the above inequality together with Lemma 2 we get

$$\mu' + \frac{3(1-q)}{2q \cdot (1+|\mu|)^2} \cdot \left( \frac{b \cdot (2+2p)}{p} - 2\mu b \right) \geq \mu + \frac{1-q}{q \cdot (1+|\mu|)^2} \cdot \nu$$

which implies that

$$b \cdot \left( 1 + \frac{3(1-q)}{2q \cdot (1+|\mu|)^2} \cdot \left( \frac{2+2p}{p} - 2\mu \right) \right) \geq \frac{1-q}{q \cdot (1+|\mu|)^2} \cdot \nu$$

and finally we get:

$$\mathbb{E}[f(\bar{Y})] - \mathbb{E}[f(\bar{X})] = b \geq \frac{p}{3+p \cdot \left( \frac{q}{1-q} \cdot (1+|\mu|)^2 + 3-3\mu \right)} \cdot \nu.$$

$\square$

## 5.2 Biasing Boolean Functions: Proving Theorem 3.11

In this section we will prove Theorem 3.11. Our proof follows ideas from the greedy attack of [1] and extend them to the blockwise setting. Despite using a similar attack, the analysis of our attack leads to weaker bounds because of being in the blockwise setting. In fact, to see why the analysis needs to be different, note that as we saw in Proposition C.3, as opposed to the binary setting achieving bias of  $+p$  in the blockwise setting for balanced functions is now *impossible* in general. In addition, making it efficient becomes further challenging in the blockwise setting. The reason is that as opposed to the binary setting, for large alphabets where we had only two choices, when the block sizes grow, we might never get a chance to sample the best possible choice for the tampered block.

In the following, we will first prove Theorem 3.11 using an inefficient (ideal greedy) tampering algorithms, and then we will show how to make it efficient with some little loss in the final bias.

### 5.2.1 Part 1: Ideal (Inefficient) Greedy Tampering

**Construction 5.12** (Ideal Greedy). Let  $\bar{X} = (X_1, \dots, X_n)$  be a joint distribution and  $f: \text{Supp}(\bar{X}) \mapsto \{+1, -1\}$ . The *ideal greedy* tampering algorithm GTam works as follows.<sup>8</sup> Given a valid prefix  $y_{\leq i-1} \in \text{ValPref}(\bar{X})$ , it holds that  $\text{GTam}(y_{\leq i-1}) = y_i$  where  $y_i$  is such that  $\hat{f}[y_{\leq i}] \geq \hat{f}[(y_1, \dots, y_{i-1}, y'_i)]$  for all valid prefixes  $(y_1, \dots, y_{i-1}, y'_i) \in \text{ValPref}(\bar{X})$ . If there are multiple such  $y_i$ , we choose the first one.

**Lemma 5.13.** *For every  $p \leq 1$ , the tampering algorithm of Construction 5.12 proves Part 1 of Theorem 3.11. Namely, it gives a robust  $p$ -tampering attack that biases  $f(\bar{X})$  up by  $\geq \frac{p \cdot (1 - \mu^2)}{2 - p \cdot (1 - \mu)}$  where  $\mu = \mathbb{E}[f(\bar{X})]$ .*

In the rest of this subsection we prove Lemma 5.13.

**Claim 5.14** (Ideal Greedy's local gains). *For any  $y_{\leq i-1} \in \text{Supp}(X_{\leq i-1})$ , it holds that*

$$\mathcal{G}_{\text{GTam}}[y_{\leq i-1}] \geq \frac{1}{2} \cdot \mathcal{A}_{X_i}[y_{\leq i-1}]$$

where GTam is the ideal greedy tampering of Construction 5.12.

*Proof.* Let us define

$$\mathcal{A}_{X_i}^+[y_{\leq i-1}] = \sum_{y_i: g[y_{\leq i}] > 0} \left[ \Pr[y_i = X_i \mid y_{\leq i-1}] \cdot |g[y_{\leq i-1}]| \right] \text{ and}$$

$$\mathcal{A}_{X_i}^-[y_{\leq i-1}] = \sum_{y_i: g[y_{\leq i}] < 0} \left[ \Pr[y_i = X_i \mid y_{\leq i-1}] \cdot |g[y_{\leq i-1}]| \right].$$

It holds that  $\mathcal{A}_{X_i}[y_{\leq i-1}] = \mathcal{A}_{X_i}^+[y_{\leq i-1}] + \mathcal{A}_{X_i}^-[y_{\leq i-1}]$ . Also, by Proposition 5.3 it holds that  $\mathcal{G}_{X_i}[y_{\leq i-1}] = \mathcal{A}_{X_i}^+[y_{\leq i-1}] - \mathcal{A}_{X_i}^-[y_{\leq i-1}] = 0$ . Therefore it holds that

$$\mathcal{A}_{X_i}[y_{\leq i-1}] = 2 \cdot \mathcal{A}_{X_i}^+[y_{\leq i-1}]. \tag{2}$$

---

<sup>8</sup>The ideal greedy could be defined for real-output functions as well, but as observed in [1], it will not lead to successful biasing attacks for general real-output functions.



On the other hand, we have:

$$\begin{aligned}
\mathcal{A}_{X_i}^+[y_{\leq i-1}] &= \sum_{y_i: g[y_{\leq i}] > 0} \Pr[y_i = X_i \mid y_{\leq i-1}] \cdot g[y_{\leq i-1}] \\
(\text{GTam maximizes } g[y_{\leq i}]) &\leq \sum_{y_i: g[y_{\leq i}] > 0} \Pr[y_i = X_i \mid y_{\leq i-1}] \cdot \mathcal{G}_{\text{GTam}}[y_{\leq i-1}] \\
&= \mathcal{G}_{\text{GTam}}[y_{\leq i-1}] \cdot \Pr_{y_i \leftarrow (Y_i | y_{\leq i-1})}[g[y_{\leq i}] > 0] \leq \mathcal{G}_{\text{GTam}}[y_{\leq i-1}].
\end{aligned}$$

Therefore by Equation (2), we get  $\mathcal{A}_{X_i}[y_{\leq i-1}]/2 \leq \mathcal{G}_{\text{GTam}}[y_{\leq i-1}]$  which proves the claim.  $\square$

**Claim 5.15.** For any  $\rho$ -tampering variation  $\bar{Y}$  of  $\bar{X}$  generated by (potentially randomized) tampering algorithm Tam, and for any  $y_{\leq i-1} \in \text{ValPref}(\bar{X})$  it holds that

$$\mathcal{A}_{Y_i}[y_{\leq i-1}] = \rho[y_{\leq i-1}] \cdot \mathcal{A}_{\text{Tam}}[y_{\leq i-1}] + (1 - \rho[y_{\leq i-1}]) \cdot \mathcal{A}_{X_i}[y_{\leq i-1}].$$

*Proof.* The proof, similarly to that of Claim 5.5, directly follows from the definition of the distribution  $Y_i$ . Namely with probability  $\rho[y_{\leq i-1}]$  we will get the outcome of  $\text{Tam}(y_{\leq i-1})$  and with probability  $1 - \rho[y_{\leq i-1}]$  we sample from the original untampered distribution  $X_i$ .  $\square$

**Claim 5.16.** Let  $\bar{X} = (X_1, \dots, X_n)$  be a joint distribution and  $f: \text{Supp}(\bar{X}) \mapsto \{+1, -1\}$  be any Boolean function with  $\mathbb{E}[f(\bar{X})] = \mu$ . Then for any distribution  $\bar{Y}$  that  $\text{Supp}(\bar{Y}) \subseteq \text{Supp}(\bar{X})$ , and  $d = \mathbb{E}[f(\bar{Y})] - \mu$  it holds that:  $\sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \mathcal{A}_{Y_i}[Y_{\leq i-1}] \geq 1 - \mu^2 - d\mu$ .

*Proof.* Firstly, we have

$$\sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \mathcal{A}_{Y_i}[Y_{\leq i-1}] = \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \mathbb{E}_{(Y_i | Y_{\leq i-1})} [|g[Y_{\leq i}]|] = \sum_{i \in [n]} \mathbb{E}_{\bar{Y}} [|g[Y_{\leq i}]|] = \mathbb{E}_{\bar{Y}} \left[ \sum_{i \in [n]} |g[Y_{\leq i}]| \right].$$

On the other hand, by the triangle inequality for every  $y \in \text{Supp}(\bar{Y})$  we have  $\sum_{i \in [n]} |g[y_{\leq i}]| \geq |\sum_{i \in [n]} g[y_{\leq i}]| = |f(y) - \mu|$ . Therefore, if  $f(y) = +1$ , then  $\sum_{i \in [n]} |g[y_{\leq i}]| \geq 1 - \mu$ , and if  $f(y) = -1$ , then  $\sum_{i \in [n]} |g[y_{\leq i}]| \geq 1 + \mu$ . Thus, we have:

$$\begin{aligned}
\sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \mathcal{A}_{Y_i}[Y_{\leq i-1}] &\geq \Pr[f(\bar{Y}) = +1] \cdot (1 - \mu) + \Pr[f(\bar{Y}) = -1] \cdot (1 + \mu) \\
&= \frac{1 + \mu + d}{2} \cdot (1 - \mu) + \frac{1 - \mu - d}{2} \cdot (1 + \mu) = 1 - \mu^2 - d\mu.
\end{aligned}$$

$\square$

Putting things together, we can conclude Lemma 5.13 as follows:

$$\begin{aligned}
& \text{(by Claim 5.16)} \quad 1 - \mu^2 - d\mu \leq \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \mathcal{A}_{Y_i}[Y_{\leq i-1}] \\
& \text{(by Claim 5.15)} \quad = \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ \rho[y_{\leq i-1}] \cdot \mathcal{A}_{\text{GTam}}[y_{\leq i-1}] \right] \\
& \quad + \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ (1 - \rho[y_{\leq i-1}]) \cdot \mathcal{A}_{X_i}[y_{\leq i-1}] \right] \\
& \text{(by Claim 5.14)} \quad \leq \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ \rho[y_{\leq i-1}] \cdot \mathcal{A}_{\text{GTam}}[y_{\leq i-1}] \right] \\
& \quad + \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ 2 \cdot (1 - \rho[y_{\leq i-1}]) \cdot \mathcal{G}_{\text{GTam}}[y_{\leq i-1}] \right] \\
& \text{(GTam always gives non-negative gain)} \quad = \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ (2 - \rho[y_{\leq i-1}]) \cdot \mathcal{G}_{\text{GTam}}[y_{\leq i-1}] \right] \\
& \text{(by Claim 5.5)} \quad = \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ \frac{2 - \rho[y_{\leq i-1}]}{\rho[y_{\leq i-1}]} \cdot \mathcal{G}_{Y_i}[y_{\leq i-1}] \right] \\
& \text{(\mathcal{G}_{Y_i}[y_{\leq i-1}] \geq 0 \text{ and } p \leq \rho[y_{\leq i-1}])} \quad \leq \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ \frac{2-p}{p} \cdot \mathcal{G}_{Y_i}[y_{\leq i-1}] \right] \\
& \text{(by Claim 5.4)} \quad = \left( \frac{2-p}{p} \right) \cdot d.
\end{aligned}$$

Therefore,  $d \geq \frac{p(1-\mu^2)}{2-p(1-\mu)}$  which finishes the proof of Lemma 5.13.

## 5.2.2 Part 2: Efficient Greedy Tampering

In this section we prove the Part 2 of Theorem 3.11. The high level idea of the proof is to try to approximate the ideal greedy algorithm by multiple samples and using Hoeffding bound.

**Construction 5.17** ( $\ell$ -Greedy tampering). Let  $\bar{X} = (X_1, \dots, X_n)$  be a joint distribution. Given a function  $f: \text{Supp}(\bar{X}) \mapsto \{+1, -1\}$  the (efficient)  $\ell$ -greedy tampering algorithm  $\text{EGTam}_\ell$  (or simply  $\text{EGTam}$  when  $\ell$  is clear from the context) works as follows. Suppose we are given a valid prefix  $y_{\leq i-1} \in \text{ValPref}(\bar{X})$ . Then:

1. Sample  $\ell$  independent instances  $y_{i,1}, \dots, y_{i,\ell}$  from  $S(y_{\leq i-1}) \equiv (X_i \mid y_{\leq i-1})$ .
2. For each  $y_{i,j}$ , sample  $\ell$  independent  $x_{i,j,1}, \dots, x_{i,j,\ell}$  from  $S(y_{\leq i-1}, y_{i,j}) \equiv (\bar{X} \mid (y_{\leq i-1}, y_{i,j}))$ .
3. Let  $\mu_{i,j} = \mathbb{E}_{k \leftarrow [\ell]} f(x_{i,j,k})$ .
4. Let  $j'$  be the smallest number in  $[\ell]$  such that  $\mu_{i,j'} \geq \mu_{i,j''}$  for any  $j' \neq j'' \in [\ell]$ .
5. Output  $y_{i,j'}$ .

**Lemma 5.18.** *There is  $\ell = \text{poly}(\frac{n}{\epsilon})$  such that for every  $p \leq 1$ , the  $\ell$ -greedy tampering algorithm of Construction 5.17 proves Part 2 of Theorem 3.11. Namely, for every  $[p, 1]$ -probability tree  $\rho$  over  $\bar{X}$ , it generates a  $p$ -tampering variation  $\bar{Y}$  of  $\bar{X}$  such that  $\mathbb{E}[f(\bar{Y})] - \mu \geq \frac{p(1-\mu^2)}{2-p(1-\mu)} - \frac{\epsilon}{1+\mu}$  where  $\mu = \mathbb{E}[f(\bar{X})]$ .*

In the rest of this section we prove Lemma 5.18.

**Definition 5.19** ( $\delta$  upper margin). We define  $\max_\delta[y_{\leq i-1}] = y'_i$  to be the “smallest”  $y'_i \in \text{Supp}(Y_i)$  such that (1)  $\Pr_{y_i \leftarrow (Y_i | y_{\leq i-1})} [\hat{f}[y_{\leq i}] \geq \hat{f}[y_{\leq i-1}, y'_i]] \geq \delta$  and (2)  $\Pr_{y_i \leftarrow (Y_i | y_{\leq i-1})} [\hat{f}[y_{\leq i}] > \hat{f}[y_{\leq i-1}, y'_i]] < \delta$ . Note that such  $y'_i$  always exists. Now we define two sets:

$$\text{Max}_{\geq \delta}[y_{\leq i-1}] = \{y \in \text{Supp}(Y_i | y_{\leq i-1}) \mid \hat{f}[y_{\leq i}] \geq \hat{f}[y_{\leq i-1}, y'_i]\}$$

$$\text{Max}_{< \delta}[y_{\leq i-1}] = \{y \in \text{Supp}(Y_i | y_{\leq i-1}) \mid \hat{f}[y_{\leq i}] > \hat{f}[y_{\leq i-1}, y'_i]\}.$$

Note that we have  $\text{Max}_{< \delta}[y_{\leq i-1}] \subset \text{Max}_{\geq \delta}[y_{\leq i-1}]$  and by the definition of  $y'_i = \max_\delta[y_{\leq i-1}]$  it holds that  $\Pr_{y_i \leftarrow (Y_i | y_{\leq i-1})} [y \in \text{Max}_{\geq \delta}[y_{\leq i-1}]] \geq \delta$  and  $\Pr_{y_i \leftarrow (Y_i | y_{\leq i-1})} [y \in \text{Max}_{< \delta}[y_{\leq i-1}]] < \delta$ .

**Definition 5.20** ( $\delta$ -greedy tampering). We call a (potentially randomized) tampering algorithm Tam a  $\delta$ -greedy algorithm if for every valid prefix  $y_{\leq i-1} \in \text{ValPref}(\bar{X})$ , and every  $y_i \leftarrow \text{Tam}(y_{\leq i-1})$  we have:

$$\hat{f}[y_{\leq i}] + \delta \geq \hat{f}[y_{\leq i-1}, \max_\delta(y_{\leq i-1})].$$

**Proposition 5.21.** For any valid prefix  $y_{\leq i-1} \in \text{ValPref}(\bar{X})$ , and  $y_i \leftarrow \text{Tam}(y_{\leq i-1})$  where Tam is a  $\delta$ -greedy tampering algorithm, we have  $g[y_{\leq i}] \geq -3 \cdot \delta$ .

*Proof.* We have

$$\begin{aligned} \hat{f}[y_{\leq i-1}] &= \sum_{y_i} \Pr[y_i = X_i \mid y_{\leq i-1}] \cdot \hat{f}[y_{\leq i}] \\ &= \sum_{y_i \in \text{Max}_{< \delta}} \Pr[y_i = X_i \mid y_{\leq i-1}] \cdot \hat{f}[y_{\leq i}] + \sum_{y_i \notin \text{Max}_{< \delta}} \Pr[y_i = X_i \mid y_{\leq i-1}] \cdot \hat{f}[y_{\leq i}] \\ &\leq \sum_{y_i \in \text{Max}_{< \delta}} \Pr[y_i = X_i \mid y_{\leq i-1}] + \sum_{y_i \notin \text{Max}_{< \delta}} \Pr[y_i = X_i \mid y_{\leq i-1}] \cdot \hat{f}[y_{\leq i-1}, \max_\delta[y_{\leq i-1}]] \\ &< \delta + (1 - \delta) \cdot \hat{f}[y_{\leq i-1}, \max_\delta[y_{\leq i-1}]] \\ &\leq \hat{f}[y_{\leq i-1}, \max_\delta[y_{\leq i-1}]] + 2 \cdot \delta. \end{aligned}$$

Therefore,

$$g[y_{\leq i}] = \hat{f}[y_{\leq i}] - \hat{f}[y_{\leq i-1}] \geq \hat{f}[y_{\leq i}] - (\hat{f}[y_{\leq i-1}, \max_\delta[y_{\leq i-1}]] + 2\delta) \geq -3 \cdot \delta.$$

□

**Claim 5.22** ( $\delta$ -greedy’s local gains). Let Tam be a  $\delta$ -greedy tampering algorithm. Then for any  $y_{\leq i-1} \in \text{Supp}(\bar{X}_{\leq i-1})$ , it holds that

$$\mathcal{G}_{\text{Tam}}[y_{\leq i-1}] \geq \frac{1}{2} \cdot \mathcal{A}_{X_i}[y_{\leq i-1}] - 7\delta.$$

*Proof.* Let us define

$$\mathcal{A}_{X_i}^*[y_{\leq i-1}] = \sum_{y_i \in \text{Max}_{< \delta}[y_{\leq i-1}]} \Pr[y_i = X_i \mid y_{\leq i}] \cdot |g[y_{\leq i}]|$$

$$\begin{aligned}\mathcal{A}_{X_i}^+[y_{\leq i-1}] &= \sum_{y_i \notin \text{Max}_{<\delta}[y_{\leq i-1}], g[y_{\leq i}] > 0} \Pr[y_i = X_i \mid y_{\leq i}] \cdot |g[y_{\leq i}]| \\ \mathcal{A}_{X_i}^-[y_{\leq i-1}] &= \sum_{y_i \notin \text{Max}_{<\delta}[y_{\leq i-1}], g[y_{\leq i}] < 0} \Pr[y_i = X_i \mid y_{\leq i}] \cdot |g[y_{\leq i}]|.\end{aligned}$$

It holds that  $\mathcal{A}_{X_i}[y_{\leq i-1}] = \mathcal{A}_{X_i}^+[y_{\leq i-1}] + \mathcal{A}_{X_i}^-[y_{\leq i-1}] + \mathcal{A}_{X_i}^*[y_{\leq i-1}]$  and at the same time (by Proposition 5.3)  $0 = \mathcal{G}_{X_i}[y_{\leq i-1}] \leq \mathcal{A}_{X_i}^+[y_{\leq i-1}] + \mathcal{A}_{X_i}^*[y_{\leq i-1}] - \mathcal{A}_{X_i}^-[y_{\leq i-1}]$ . Therefore

$$\mathcal{A}_{X_i}[y_{\leq i-1}] \leq 2 \cdot \mathcal{A}_{X_i}^+[y_{\leq i-1}] + 2 \cdot \mathcal{A}_{X_i}^*[y_{\leq i-1}]. \quad (3)$$

For  $\mathcal{A}_{X_i}^+[y_{\leq i-1}]$  we have

$$\begin{aligned}\mathcal{A}_{X_i}^+[y_{\leq i-1}] &\leq \sum_{y_i \notin \text{Max}_{<\delta}[y_{\leq i-1}], g[y_{\leq i}] > 0} \Pr[y_i = Y_i \mid y_{\leq i}] \cdot g[y_{\leq i}] \\ &\leq \sum_{y_i \notin \text{Max}_{<\delta}[y_{\leq i-1}], g[y_{\leq i}] > 0} \Pr[y_i = Y_i \mid y_{\leq i}] \cdot g[y_{\leq i-1}, \text{max}_\delta[y_{\leq i-1}]] \\ &= g[y_{\leq i-1}, \text{max}_\delta[y_{\leq i-1}]] \cdot \sum_{y_i \notin \text{Max}_{<\delta}[y_{\leq i-1}], g[y_{\leq i}] > 0} \Pr[y_i = Y_i \mid y_{\leq i-1}] \\ &\leq (\mathcal{G}_{\text{Tam}}[y_{\leq i-1}] + \delta) \cdot \sum_{y_i \notin \text{Max}_{<\delta}[y_{\leq i-1}], g[y_{\leq i}] > 0} \Pr[y_i = Y_i \mid y_{\leq i-1}] \\ &\leq |\mathcal{G}_{\text{Tam}}[y_{\leq i-1}] + \delta|.\end{aligned}$$

Moreover, for  $\mathcal{A}_{X_i}^*[y_{\leq i-1}]$  we have

$$\begin{aligned}\mathcal{A}_{X_i}^*[y_{\leq i-1}] &\leq \sum_{y_i \in \text{Max}_{<\delta}[y_{\leq i-1}], g[y_{\leq i}] > 0} \Pr[y_i = Y_i \mid y_{\leq i}] \cdot g[y_{\leq i}] \\ (g[y_{\leq i}] \text{ is at most } 2) &\leq \sum_{y_i \in \text{Max}_{<\delta}[y_{\leq i-1}], g[y_{\leq i}] > 0} \Pr[y_i = Y_i \mid y_{\leq i}] \cdot 2 \\ &\leq 2 \cdot \Pr_{y_i \leftarrow (Y_i | y_{\leq i})} [g[y_{\leq i}] > 0 \wedge y_i \in \text{Max}_{<\delta}[y_{\leq i-1}]] < 2 \cdot \delta.\end{aligned}$$

Therefore by Equation (3), we have:

$$\begin{aligned}\mathcal{A}_{X_i}[y_{\leq i-1}] &\leq 2 \cdot |\mathcal{G}_{\text{Tam}}[y_{\leq i-1}] + \delta| + 4 \cdot \delta \\ (\text{by proposition 5.21}) &\leq 2 \cdot (\mathcal{G}_{\text{Tam}}[y_{\leq i-1}] + 5 \cdot \delta) + 4 \cdot \delta \\ &\leq 2 \cdot \mathcal{G}_{\text{Tam}}[y_{\leq i-1}] + 14 \cdot \delta\end{aligned}$$

which proves the claim.  $\square$

**Claim 5.23** ( $\delta$ -greedy's global gain). *Any  $\delta$ -greedy tampering algorithm Tam is a robust  $p$ -tampering attack for biasing  $f(\bar{X})$  by  $\frac{p \cdot (1 - \mu^2)}{2 - p \cdot (1 - \mu)} - \frac{17 \cdot \delta \cdot n}{1 + \mu}$  where  $\mu = \mathbb{E}[f(\bar{X})]$ .*

*Proof.* First we give a lower bound on  $d = \mathbb{E}[f(\bar{Y})] - \mu$  for  $\delta$ -greedy tampering algorithms.

$$\begin{aligned}
& \text{(by Claim 5.16)} \quad 1 - \mu^2 - d\mu \leq \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \mathcal{A}_{Y_i}[\bar{Y}_{\leq i-1}] \\
& \text{(by Claim 5.15)} \quad = \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ \rho[y_{\leq i-1}] \cdot \mathcal{A}_{\text{Tam}}[y_{\leq i-1}] + (1 - \rho[y_{\leq i-1}]) \cdot \mathcal{A}_{X_i}[y_{\leq i-1}] \right] \\
& \text{(by Claim 5.22)} \quad \leq \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ \rho[y_{\leq i-1}] \cdot \mathcal{A}_{\text{Tam}}[y_{\leq i-1}] \right] \\
& \quad \quad \quad + \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ 2 \cdot (1 - \rho[y_{\leq i-1}]) \cdot (\mathcal{G}_{\text{Tam}}[y_{\leq i-1}] + 7\delta) \right] \\
& \text{(by Prop. 5.21)} \quad \leq \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ \rho[y_{\leq i-1}] \cdot (\mathcal{G}_{\text{Tam}}[y_{\leq i-1}] + 6 \cdot \delta) \right] \\
& \quad \quad \quad + \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ 2 \cdot (1 - \rho[y_{\leq i-1}]) \cdot (\mathcal{G}_{\text{Tam}}[y_{\leq i-1}] + 7\delta) \right] \\
& \quad \quad \quad = \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ (2 - \rho[y_{\leq i-1}]) \cdot \mathcal{G}_{\text{Tam}}[y_{\leq i-1}] + (14 - 8 \cdot \rho[y_{\leq i-1}]) \cdot \delta \right] \\
& \text{(by Claim 5.5)} \quad = \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ \frac{2 - \rho[y_{\leq i-1}]}{\rho[y_{\leq i-1}]} \cdot \mathcal{G}_{Y_i}[y_{\leq i-1}] + (14 - 8 \cdot \rho[y_{\leq i-1}]) \cdot \delta \right] \\
& \quad \quad \quad \leq \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ \frac{2 - \rho[y_{\leq i-1}]}{\rho[y_{\leq i-1}]} \cdot (\mathcal{G}_{Y_i}[y_{\leq i-1}] + 3 \cdot \delta) \right] \\
& \quad \quad \quad \quad - \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ 3 \cdot \delta + (14 - 8 \cdot \rho[y_{\leq i-1}]) \cdot \delta \right] \\
& \text{(Prop. 5.21 and } \rho[y_{\leq i-1}] \geq p) \quad \leq \sum_{i \in [n]} \mathbb{E}_{Y_{\leq i-1}} \left[ \frac{2-p}{p} \cdot (\mathcal{G}_{Y_i}[y_{\leq i-1}] + 3 \cdot \delta) + (11 - 8 \cdot \rho[y_{\leq i-1}]) \cdot \delta \right] \\
& \text{(by Claim 5.4)} \quad \leq \left( \frac{2-p}{p} \right) \cdot d + \frac{17 \cdot \delta \cdot n}{p}.
\end{aligned}$$

Therefore we have

$$d \geq \frac{p \cdot (1 - \mu^2)}{2 - p \cdot (1 - \mu)} - \frac{17 \cdot \delta \cdot n}{2 - p \cdot (1 - \mu)} \geq \frac{p \cdot (1 - \mu^2)}{2 - p \cdot (1 - \mu)} - \frac{17 \cdot \delta \cdot n}{(1 + \mu)}.$$

□

*Proof of Lemma 5.18.* First, fix  $i \in [n]$ . With probability  $1 - (1 - \delta)^\ell$ , there exist  $j'' \in [\ell]$  such that  $y_{i,j''} \in \text{Max}_{\geq \delta}[y_{\leq i-1}]$ . Also, for every fixed  $j \in [\ell]$ , by Hoeffding bound,

$$\Pr \left[ |u_{i,j} - \hat{f}[y_{\leq i-1}, y_{i,j}]| \geq \frac{\delta}{2} \right] \leq 2 \cdot e^{-\ell \cdot \frac{\delta^2}{8}}.$$

By a union bound, with probability at least  $1 - (1 - \delta)^\ell - (\ell)(2 \cdot e^{-\ell \cdot \frac{\delta^2}{8}}) \geq 1 - 3\ell e^{-\ell \cdot \frac{\delta^2}{8}}$  we will have both (1)  $j'' \in [\ell]$  such that  $y_{i,j''} \in \text{Max}_{\geq \delta}[y_{\leq i-1}]$  and (2) all the approximations are  $\delta/2$  correct. If both of these

conditions hold, the  $\ell$ -greedy will be guaranteed to pick some  $y_{i,j'}$  that satisfies the  $\delta$ -greedy definition for the fixed level  $i \in [n]$ . By picking  $\ell = \text{poly}(1/\delta)$  (e.g.,  $\ell \gg 1/\delta^3$  suffices) the probability of failing in level  $i$  will be at most  $3\ell e^{-\ell \cdot \frac{\delta^2}{8}} < \delta/2 < \delta/(1+\mu)$ , and by a union bound over  $i \in [n]$ , the total probability of failing at *any* level will be at most  $n\delta/(1+\mu)$ . So, for sufficiently large  $\ell \geq \text{poly}(1/\delta)$ , the  $\ell$ -greedy is  $\delta$ -greedy with probability at least  $1 - \frac{n \cdot \delta}{1+\mu}$ , and so the total amount of bias for  $\ell$ -greedy will be least

$$\left(1 - \frac{n \cdot \delta}{1+\mu}\right) \cdot \left(\frac{p \cdot (1-\mu^2)}{2-p \cdot (1-\mu)} - \frac{17 \cdot \delta \cdot n}{1+\mu}\right) - 2 \cdot \frac{n \cdot \delta}{1+\mu} \geq \frac{p \cdot (1-\mu^2)}{2-p \cdot (1-\mu)} - \frac{20 \cdot \delta \cdot n}{1+\mu}.$$

Now setting  $\delta = \frac{\varepsilon}{20n}$ , we get  $\ell = \text{poly}(n/\varepsilon)$ , and that finishes the proof.  $\square$

## 6 Open Questions

We conclude by describing some open questions and interesting directions for future research.

**Power of  $k$ -sampling attacks for small  $k$ .** A natural yet more general class of attacks that include  $k$ -resetting attacks at special case is the class of  $k+1$  sampling attacks in which the tampering algorithm first gets  $k+1$  samples from the distribution of the  $i$ 'th tampered block and then it chooses one of these samples (perhaps by calls to the online sampler and the function  $f$ ). Our  $\ell$ -greedy algorithm is indeed an  $\ell$  sampling attack but to get good bias, it needs to use many  $\ell = \text{poly}(n/\varepsilon)$  samples. What is the power of  $\ell$ -sampling attacks in general, when  $\ell$  is small, e.g. constant?

**Power of ‘very efficient’ viruses.** What is the power of tampering attacks whose computational resources is not sufficient for sampling the next block or even computing  $f$ ? Such tampering algorithms are natural for cryptographic attacks where computing  $f$  is heavy and the virus might prefer to use very limited resources not to be detected by the system. Our efficient tampering attacks of Theorems 3.9 and 3.11 both need to run the online sampler as well as the function  $f$ . It remains an interesting future direction to study the power of limited tampering attacks whose decisions are more ‘local’ and cannot be based on sampling the blocks from the original distribution or computing  $f$ . We conjecture that such efficient viruses that cannot depend on  $f$  or the distribution  $\bar{X}$  are not powerful to achieve constant bias  $\Omega(p)$ . However, it is interesting to find out what is the *minimum* number of calls needed to  $f$  or the sampler for getting bias  $\Omega(p)$ .

**Optimal constant.** In Proposition C.3 we show that blockwise  $p$ -tampering attacks cannot achieve bias  $c \cdot p$  for all  $p$  if  $c > (1+\mu) \cdot \ln(\frac{2}{1+\mu})$ . This leaves open finding the optimal constant  $c$  for which  $p$ -tampering attacks can always get  $\geq c \cdot p$  bias. For balanced Boolean functions we already know  $0.5 \leq c \leq \ln(2) < 0.7$ .

**Biasing up vs. biasing either way.** Our Theorems 3.9 and 3.11 always bias the function towards  $+1$ . Inspired by models of attacks against coin-tossing protocols [8, 14, 15, 17, 29, 34] one can ask the following questions. What is the power of  $p$ -tampering biasing attacks whose goal is to *either* bias the average of the function up *or* bias it down? Some of the applications of our biasing attacks (e.g., against learners) need to bias the function always in a fixed direction to increase the ‘error’, but other attacks (e.g., against extractors) could achieve their goal by biasing the function in either direction. Our Proposition C.3 does not apply to this setting, as it only limits the power of biasing attacks who always want to bias the function upwards.

**Acknowledgement.** We thank Dimitrios Diochnos, Yevgeniy Dodis, and Yanjun Qi for useful discussions.

## References

- [1] Per Austrin, Kai-Min Chung, Mohammad Mahmoody, Rafael Pass, and Karn Seth. On the impossibility of cryptography with tamperable randomness. In *International Cryptology Conference*, pages 462–479. Springer, 2014. [1](#), [3](#), [4](#), [5](#), [8](#), [9](#), [14](#), [15](#), [16](#), [21](#), [32](#), [42](#), [43](#), [45](#), [46](#)
- [2] Pranjal Awasthi, Maria Florina Balcan, and Philip M Long. The power of localization for efficiently learning linear separators with noise. In *Proceedings of the 46th Annual ACM Symposium on Theory of Computing*, pages 449–458. ACM, 2014. [6](#), [7](#)
- [3] Yossi Azar, Andrei Z Broder, Anna R Karlin, Nathan Linial, and Steven Phillips. Biased random walks. In *Proceedings of the twenty-fourth annual ACM symposium on Theory of computing*, pages 1–9. ACM, 1992. [4](#)
- [4] Boaz Barak and Shai Halevi. A model and architecture for pseudo-random generation with applications to/dev/random. In *Proceedings of the 12th ACM conference on Computer and communications security*, pages 203–212. ACM, 2005. [3](#)
- [5] Marco Barreno, Blaine Nelson, Anthony D Joseph, and JD Tygar. The security of machine learning. *Machine Learning*, 81(2):121–148, 2010. [6](#)
- [6] Salman Beigi, Omid Etesami, and Amin Gohari. Deterministic randomness extraction from generalized and distributed santha–vazirani sources. *SIAM Journal on Computing*, 46(1):1–36, 2017. [5](#), [6](#), [11](#), [17](#), [45](#)
- [7] Mihir Bellare, Kenneth G Paterson, and Phillip Rogaway. Security of symmetric encryption against mass surveillance. In *Advances in Cryptology–CRYPTO 2014*, pages 1–19. Springer, 2014. [9](#)
- [8] Itay Berman, Iftach Haitner, and Aris Tentes. Coin flipping of any constant bias implies one-way functions. In *Proceedings of the 46th Annual ACM Symposium on Theory of Computing*, pages 398–407. ACM, 2014. [9](#), [38](#)
- [9] Boneh, DeMillo, and Lipton. On the importance of checking cryptographic protocols for faults. In *EUROCRYPT: Advances in Cryptology: Proceedings of EUROCRYPT*, 1997. [9](#)
- [10] Nader H Bshouty, Nadav Eiron, and Eyal Kushilevitz. Pac learning with nasty noise. *Theoretical Computer Science*, 288(2):255–275, 2002. [6](#), [7](#)
- [11] Nishanth Chandran, Vipul Goyal, Pratyay Mukherjee, Omkant Pandey, and Jalaj Upadhyay. Block-wise non-malleable codes. In *LIPICs-Leibniz International Proceedings in Informatics*, volume 55. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2016. [9](#)
- [12] Benny Chor and Oded Goldreich. Unbiased bits from sources of weak randomness and probabilistic communication complexity. In *Proc. 26th FOCS*, pages 429–442. IEEE, 1985. [5](#)
- [13] Benny Chor and Oded Goldreich. Unbiased bits from sources of weak randomness and probabilistic communication complexity. *SIAM Journal on Computing*, 17(2):230–261, 1988. [6](#), [11](#)
- [14] Richard Cleve. Limits on the security of coin flips when half the processors are faulty. In *Proceedings of the eighteenth annual ACM symposium on Theory of computing*, pages 364–369. ACM, 1986. [38](#)

- [15] Richard Cleve and Russell Impagliazzo. Martingales, collective coin flipping and discrete control processes. *In other words*, 1:5, 1993. 9, 38
- [16] Henry Corrigan-Gibbs and Suman Jana. Recommendations for randomness in the operating system, or how to keep evil children out of your pool and other random facts. In *HotOS*, 2015. 3
- [17] Dana Dachman-Soled, Yehuda Lindell, Mohammad Mahmoody, and Tal Malkin. On the black-box complexity of optimally-fair coin tossing. In *Theory of Cryptography Conference*, pages 450–467. Springer, 2011. 9, 38
- [18] Ivan Damgård, Sebastian Faust, Pratyay Mukherjee, and Daniele Venturi. Tamper resilient cryptography without self-destruct. Cryptology ePrint Archive, Report 2013/124, 2013. <http://eprint.iacr.org/2013/124>. 4
- [19] Dodis, Ong, Prabhakaran, and Sahai. On the (im)possibility of cryptography with imperfect randomness. In *FOCS: IEEE Symposium on Foundations of Computer Science (FOCS)*, 2004. 1, 3, 5, 6, 45
- [20] Yevgeniy Dodis. New imperfect random source with applications to coin-flipping. *Automata, Languages and Programming*, pages 297–309, 2001. 9
- [21] Yevgeniy Dodis, David Pointcheval, Sylvain Ruhault, Damien Vergniaud, and Daniel Wichs. Security analysis of pseudo-random number generators with input:/dev/random is not robust. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, pages 647–658. ACM, 2013. 3
- [22] Yevgeniy Dodis and Yanqing Yao. Privacy with imperfect randomness. In *Annual Cryptology Conference*, pages 463–482. Springer, 2015. 5, 6
- [23] Cynthia Dwork, Guy N Rothblum, and Salil Vadhan. Boosting and differential privacy. In *Foundations of Computer Science (FOCS), 2010 51st Annual IEEE Symposium on*, pages 51–60. IEEE, 2010. 6, 10, 11, 16, 17
- [24] Stefan Dziembowski, Sebastian Faust, and François-Xavier Standaert. Private circuits III: Hardware trojan-resilience via testing amplification. In Edgar R. Weippl, Stefan Katzenbeisser, Christopher Kruegel, Andrew C. Myers, and Shai Halevi, editors, *ACM CCS 16: 23rd Conference on Computer and Communications Security*, pages 142–153, Vienna, Austria, October 24–28, 2016. ACM Press. 4
- [25] Stefan Dziembowski, Krzysztof Pietrzak, and Daniel Wichs. Non-malleable codes. In Andrew Chi-Chih Yao, editor, *ICS*, pages 434–452. Tsinghua University Press, 2010. 9
- [26] Rosario Gennaro, Anna Lysyanskaya, Tal Malkin, Silvio Micali, and Tal Rabin. Algorithmic tamper-proof (ATP) security: Theoretical foundations for security against hardware tampering. In Moni Naor, editor, *TCC 2004: 1st Theory of Cryptography Conference*, volume 2951 of *Lecture Notes in Computer Science*, pages 258–277, Cambridge, MA, USA, February 19–21, 2004. Springer, Heidelberg, Germany. 4, 9
- [27] Shafi Goldwasser, Yael Tauman Kalai, and Sunoo Park. Adaptively secure coin-flipping, revisited. In *International Colloquium on Automata, Languages, and Programming*, pages 663–674. Springer, 2015. 9



- [28] Zvi Gutterman, Benny Pinkas, and Tzachy Reinman. Analysis of the linux random number generator. In *Security and Privacy, 2006 IEEE Symposium on*, pages 15–pp. IEEE, 2006. [3](#)
- [29] Iftach Haitner and Eran Omri. Coin flipping with constant bias implies one-way functions. *SIAM Journal on Computing*, 43(2):389–409, 2014. [9](#), [38](#)
- [30] Nadia Heninger, Zakir Durumeric, Eric Wustrow, and J Alex Halderman. Mining your ps and qs: Detection of widespread weak keys in network devices. In *USENIX Security Symposium*, volume 8, 2012. [3](#)
- [31] Michael Kearns and Ming Li. Learning in the presence of malicious errors. *SIAM Journal on Computing*, 22(4):807–837, 1993. [6](#), [7](#)
- [32] Aggelos Kiayias and Yiannis Tselekounis. Tamper resilient circuits: The adversary at the gates. In Kazue Sako and Palash Sarkar, editors, *Advances in Cryptology – ASIACRYPT 2013, Part II*, volume 8270 of *Lecture Notes in Computer Science*, pages 161–180, Bangalore, India, December 1–5, 2013. Springer, Heidelberg, Germany. [4](#)
- [33] David Lichtenstein, Nathan Linial, and Michael Saks. Some extremal problems arising from discrete control processes. *Combinatorica*, 9(3):269–287, 1989. [9](#)
- [34] Hemanta K Maji, Manoj Prabhakaran, and Amit Sahai. On the computational complexity of coin flipping. In *Foundations of Computer Science (FOCS), 2010 51st Annual IEEE Symposium on*, pages 613–622. IEEE, 2010. [9](#), [38](#)
- [35] Nicolas Papernot, Patrick McDaniel, Arunesh Sinha, and Michael Wellman. Towards the science of security and privacy in machine learning. *arXiv preprint arXiv:1611.03814*, 2016. [6](#)
- [36] Omer Reingold, Salil Vadhan, and Avi Wigderson. A note on extracting randomness from santha-vazirani sources. *Unpublished manuscript*, 2004. [1](#), [5](#), [6](#), [45](#)
- [37] Benjamin IP Rubinstein, Blaine Nelson, Ling Huang, Anthony D Joseph, Shing-hon Lau, Satish Rao, Nina Taft, and JD Tygar. Antidote: understanding and defending against poisoning of anomaly detectors. In *Proceedings of the 9th ACM SIGCOMM conference on Internet measurement conference*, pages 1–14. ACM, 2009. [6](#)
- [38] Benjamin IP Rubinstein, Blaine Nelson, Ling Huang, Anthony D Joseph, Shing-hon Lau, Satish Rao, Nina Taft, and JD Tygar. Stealthy poisoning attacks on pca-based anomaly detectors. *ACM SIGMETRICS Performance Evaluation Review*, 37(2):73–74, 2009. [6](#)
- [39] Miklos Santha and Umesh V. Vazirani. Generating quasi-random sequences from semi-random sources. *J. Comput. Syst. Sci.*, 33(1):75–87, 1986. [3](#), [5](#), [6](#), [11](#), [46](#)
- [40] Ronen Shaltiel. Recent developments in explicit constructions of extractors. *Bulletin of the EATCS*, 77(67-95):10, 2002. [43](#)
- [41] Shiqi Shen, Shruti Tople, and Prateek Saxena. A uror: defending against poisoning attacks in collaborative deep learning systems. In *Proceedings of the 32nd Annual Conference on Computer Security Applications*, pages 508–519. ACM, 2016. [6](#), [7](#)

- [42] Leslie G Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984. 18
- [43] Leslie G Valiant. Learning disjunction of conjunctions. In *IJCAI*, pages 560–566, 1985. 6
- [44] John Von Neumann. 13. various techniques used in connection with random digits. *Appl. Math Ser*, 12:36–38, 1951. 5
- [45] Huang Xiao, Battista Biggio, Gavin Brown, Giorgio Fumera, Claudia Eckert, and Fabio Roli. Is feature selection secure against training data poisoning? In *ICML*, pages 1689–1698, 2015. 6, 7

## A Blockwise $p$ -Tampering Attacks on Primitives: Case of Encryption

In this section we give a formal definition of what it means to break the security of public-key encryption through a blockwise  $p$ -tampering attack. The definition is similar to that of [1] while here we require the attacking virus to succeed for *any* partitioning of the incoming randomness and shall do it in a ‘robust’ way; namely, by only assuming that the tampering probabilities are  $\geq p$  (rather than exactly  $p$ ). The partitioning of the randomness into blocks could be due to the fact that randomness is being generated in ‘chunks’ upon request, or it might be due to some larger alphabets used by the machine to represent data. Our results directly extends to more primitives studied in [1] by plugging in our blockwise tampering attacks. Thus, we only formalize the blockwise tampering attack for the case of public-key encryption, and we refer the reader to [1] for other primitives such as private-key encryption, commitments, secure computation, zero-knowledge proofs, etc.

**Definition A.1** (Partitioning of strings). For a sequence of (possibly empty) strings  $(r_1, \dots, r_m)$  we denote their concatenation by  $(r_1 | \dots | r_m)$  where the lengths of  $r_i$ ’s is not necessarily clear anymore.<sup>9</sup> We call the sequence  $(s_1, \dots, s_n)$  a partitioning of  $[N]$  if  $\sum_i s_i = N$  and  $s_i \geq 0$  for all  $i$ . For a string  $R$  of length  $N$  and a partitioning  $S = (s_1, \dots, s_n)$  of  $[N]$  we call  $(r_1, \dots, r_m)$  the partitioning of  $R$  according to  $S$  if  $|r_i| = s_i$  and  $R$  is equal to the concatenation  $R = (r_1 | \dots | r_n)$ .

**Interpretation.** The partitioning  $(r_1, \dots, r_m)$  of a long random seed  $r$  could be as a result of various scenarios. It could simply refer to a larger alphabet size being used in the system to represent blocks of randomness. More naturally, each block could refer to what the ‘system’ generates upon the request of a randomized algorithm (e.g., encryption).

As in [1], we use the following definition which guarantees a *weak* form of security against blockwise  $p$ -tampering attacks by restricting the adversary to only tamper with the randomness of the encryption (and not the key generation). However, since our results are *negative* by showing the power of attackers, such restriction only makes our results stronger.

**Definition A.2** (Blockwise  $p$ -tampering attacks on encryption). Suppose  $\mathcal{P} = (\text{Gen}, \text{Enc}, \text{Dec})$  is a public key encryption scheme in which  $\text{Enc}(\kappa, x)$  uses randomness  $r_E$  of length  $N = \text{poly}(\kappa, |x|)$  where  $\kappa$  is the security parameter. We say that the adversary  $\mathbf{Adv}$  robustly  $\alpha$ -breaks the CPA security of  $\mathcal{P}$  through blockwise  $p$ -tampering if  $\mathbf{Adv}$  can distinguish  $b = 0$  from  $b = 1$  in the experiment below by an advantage of  $\geq \alpha$  for *any* partitioning  $s_1, \dots, s_n$  of  $[N]$ .

1. A pair of keys  $(\text{sk}, \text{pk}) \leftarrow \text{Gen}(1^\kappa)$  are generated and  $\mathbf{Adv}$  receives the public-key  $\text{pk}$ .

<sup>9</sup>This is in contrast with the notation  $(r_1, \dots, r_m)$  in which the encoding will allow separating  $r_i$  from each other.

2. **Adv** generates a (supposedly  $p$ -tampering) circuit  $T$  and two of messages  $\{x_0, x_1\}$  of equal lengths.
3. The randomness  $r_E$  gets generated in blocks of lengths  $s_1, \dots, s_n$  and the tampering algorithm  $T$  will perform a blockwise  $p$ -tampering attack against the corresponding partitioning  $(r_1, \dots, r_n)$  of the randomness  $r_E$ . In particular, each new block  $r_i$  is first sampled from  $U_{s_i}$ , and then with (independent) probability  $p$ ,  $T$  is allowed to substitute  $r_i$  with something of the same length  $s_i$ . Let  $r'_E$  be the final (possibly tampered) randomness.
4. The challenger gets  $x_b \in \{x_0, x_1\}$  and sends  $c = \text{Enc}_{\text{pk}}(x_b; r'_E)$  to **Adv**.
5. **Adv** receives  $c$ , outputs a (hopefully distinguishing) bit  $b'$ .

**Theorem A.3** (Blockwise  $p$ -tampering attacks on encryption). *For any (correct) public-key encryption  $\mathcal{P} = (\text{Gen}, \text{Enc}, \text{Dec})$  scheme there is a polynomial time adversary **Adv** who robustly  $\Omega(p)$ -breaks the CPA security of  $\mathcal{P}$  through a blockwise  $p$ -tampering.*

*Proof Sketch for Theorem A.3.* The proof follows the footsteps of [1] closely, while the only difference is that this time we use our blockwise  $p$ -tampering attacker of Theorem 3.9 for real-output functions.

1. The adversary samples its own randomness  $r$  for a *strong* seeded extractor  $\text{Ext}_r(\cdot)$  [40] that extracts a bit with bias  $\leq o(p)$  from sources with min-entropy  $\omega(\log \kappa)$ . It also trivially picks  $x_0 = 0, x_1 = 1$ .
2. Given the public-key  $\text{pk}$ , the adversary defines the functions:  $f_0(\cdot), f_1(\cdot)$  where  $f_b(\cdot)$  takes as input  $r_E$  (i.e.,  $\text{Enc}$ 's randomness) and outputs:  $\text{Ext}_r(\text{Enc}_{\text{pk}}(x_b; r_E))$ . The adversary finally lets  $f(r_E) = f_1(r_E) - f_0(r_E) \in \{-1, 0, 1\}$ .
3. The adversary generates a tampering circuit  $\text{Tam}$  of Theorem 3.9 such that biases the output of the function  $f$  towards  $+1$  by at least  $\Omega(p)$  through a blockwise  $p$  tampering attack. The only point is that the tampering circuit  $T$  does not know the length of the blocks ahead of the time. But, since  $r_i$  is actually *given* to  $T$  before  $T$  substitutes it with possible tampered  $r'_i$ ,  $T$  would discover the length  $s_i$  at that moment. Also, note that  $\text{Tam}$  of Theorem 3.9 needs to call  $f$  internally, but that is not an issue since  $f$  is known to the adversary and can generate  $T$  accordingly.
4. After receiving the actual ciphertext  $c$  (encrypted with the *tampered* randomness  $r'_E$ ), the adversary applies the function to get  $b' = \text{Ext}_r(c)$  and outputs  $b'$ .

**The analysis.** We refer the reader to [1] for the actual proof (while one has to use our blockwise tampering attacker instead), roughly speaking, the reason that the above attack succeeds is the following.

- Because of the CPA security of the original scheme, the ciphertext will have  $\omega(\log \kappa)$  bits of min-entropy which guarantees the strong extractor will indeed sample an almost unbiased bit in  $\{0, 1\}$ .
- Therefore the difference function  $f(r_E) = f_1(r_E) - f_0(r_E)$  will also be an almost unbiased function with range  $\{-1, 0, +1\} \subset [-1, +1]$ .
- The *correctness* of the original PKE scheme, implies that  $\text{Var}[f(r_E)] \geq \Omega(1)$ .
- By  $\text{Var}[f(r_E)] \geq \Omega(1)$ , due to the properties of the  $p$ -tampering attacker of Theorem 1.2 we conclude that  $f(r_E)$  will be biased towards  $+1$  by  $\Omega(p)$ , and because the original distribution had bias  $\mathbb{E}[f(r_E)] \leq o(p)$ , we would now get  $\mathbb{E}[f(r'_E)] \geq \Omega(p) - o(p) \geq \Omega(p)$ .

- Finally note that  $\mathbb{E}[f(r'_E)] \geq \Omega(p)$  is equivalent to  $\Pr[f(r'_E) = +1] \geq \Pr[f(r'_E) = -1] + \Omega(p)$ , which in turn is equivalent to  $\Pr_{r'_E}[f_1(r'_E) = 1 \wedge f_0(r'_E) = 0] - \Pr_{r'_E}[f_1(r'_E) = 0 \wedge f_0(r'_E) = 1] \geq \Omega(p)$ . Therefore,  $\Pr_{r'_E}[f_1(r'_E) = 1] \geq \Pr_{r'_E}[f_0(r'_E) = 1] + \Omega(p)$  and so the output  $b'$  will in fact distinguish  $b = 0$  from  $b = 1$  experiments with advantage  $\Omega(p) \geq 1/\text{poly}(\kappa)$ .

□

## B Reducing Blockwise Tampering to Bitwise for Uniform Distributions

**Notation.** For any partitioning  $S = [s_1, \dots, s_n]$  of  $[N]$  (according to Definition A.1) we use  $U_S$  to denote the product distribution  $U_{s_1} \times \dots \times U_{s_n}$ . (Note that  $U_S$  has  $n$  blocks, so if  $n > 1$  then  $\langle U_S \rangle \not\equiv U_N$ , and if  $n < N$  then  $\langle U_S \rangle \not\equiv U_1^N$ .) Let  $S = (s_1, \dots, s_n)$  be a partitioning of  $[N]$  and  $\bar{X} = (X_1, \dots, X_n)$  be a joint distribution such that for each  $i \in [n]$ ,  $X_i$  is a distribution over  $\{0, 1\}^{s_i}$ . We use  $\langle \bar{X} \rangle$  to denote a distribution over  $\{0, 1\}^N$  which is identical to  $\bar{X}$ , but *without* the separators between blocks; namely, they are concatenated into one block. For example, if  $S = [s_1, \dots, s_n]$  is a partitioning of  $[N]$ , then  $\langle U_S \rangle \equiv U_N$ .

**Theorem B.1** (Reducing blockwise tampering to bitwise tampering for uniform distributions). *Let  $\bar{Y}$  be a  $p$ -tampering variation of  $\bar{X} = U_1^N$  and  $S = (s_1, \dots, s_n)$  be a partitioning of  $[N]$ . If  $1 - \tilde{p} \leq (1 - p)^{\max(s_i)}$ , then there exists a  $\tilde{p}$ -tampering variation  $\tilde{Y}$  of  $U_S$  such that  $\langle \bar{Y} \rangle \equiv \langle \tilde{Y} \rangle$ , i.e., without the separators, they are the same distributions. Moreover, if  $\bar{Y}$  can be generated by efficient  $p$ -tampering Tam from  $\bar{X}$ ,  $\tilde{Y}$  can be generated by an efficient  $\tilde{p}$ -tampering Tam (who makes oracle calls to Tam) from  $U_S$ .*

*Proof.* We define  $B_p^k$  to be the product distribution of  $k$  independent Bernoulli variables that take 1 with probability  $p$ . We define an event  $E$  over a bit-string of size  $k$  as follows.  $E(b_1, \dots, b_k) = 1$  iff, at least one of  $b_i$ 's is 1. Now we define the distribution  $\hat{B}_p^k \equiv (B_p^k | E)$ . Let Tam be the tampering algorithm for  $U_1^N$  that generates  $\bar{Y}$  and let  $s = \max(s_i)$ . Now, we build a new tampering algorithm Tam for  $U_S$ . Given a valid prefix  $\tilde{y}_{\leq i-1} \in \text{ValPref}(U_S)$ , Tam does the following:

- with probability  $\alpha = ((1 - p)^{s_i} - (1 - \tilde{p})) / \tilde{p}$ 
  1. let  $(b_1, \dots, b_{s_i}) = (0, \dots, 0)$
  2. sample  $\tilde{y}_i \leftarrow U_{s_i}$  and output  $\tilde{y}_i$
- and with probability  $1 - \alpha$ 
  1. sample a bit-string  $(b_1, \dots, b_{s_i}) \leftarrow \hat{B}_p^{s_i}$
  2. for  $j \in [s_i]$ 
    - (a) if  $b_j = 0$  then  $a_j \leftarrow U_1$
    - (b) else  $a_j \leftarrow \text{Tam}([\tilde{y}_{\leq i-1}, a_{\leq j-1}])$
  3. output  $\tilde{y}_i = a$

Let  $\tilde{Y}$  be the  $\tilde{p}$ -tampering variation of  $U_S$  generated through Tam. Now we show that  $\langle \tilde{Y} \rangle \equiv \langle Y \rangle$ . Namely, we argue that Tam is perfectly simulating the  $p$ -tampering setting for Tam on each bit, which is equivalent to saying that the probability that Tam is used is exactly  $p$  for every bit *independent* of all the other bits. Note that Tam is used to tamper the  $j$ 'th bit, if and only if  $b_j = 1$ . When we use Tam in  $\tilde{p}$  tampering setting, with probability  $1 - \tilde{p} + \tilde{p} \cdot \frac{(1-p)^{s_i} - (1-\tilde{p})}{\tilde{p}} = (1 - p)^{s_i} = \Pr[\bar{E}]$  the bit-string  $(b_1, \dots, b_{s_i})$  is sampled from the distribution  $B_p^{s_i} | \bar{E}$ , and with probability  $1 - (1 - p)^{s_i} = \Pr[E]$  it is sampled from distribution  $B_p^{s_i} | E$ . Therefore, at the end  $b_1, \dots, b_s$  is indeed sampled from  $B_p^{s_i}$ , which means the tampering happens on each bit independently with probability  $p$ .

□

## C Power and Limitation of Inefficient $p$ -Tampering Attacks

First, using techniques and arguments from [6, 19, 36] we observe that the bound proved in Theorem 3.11 could be achieved even for (bounded) *real-output* functions, however, unfortunately we could only get this bias for through a computationally unbounded  $p$ -tampering attack. Due to the inefficiency of our attacks we find it more natural to explain our results as proving the existence of certain ‘tampering variations’ of the original distributions.

**Theorem C.1** (Inefficient blockwise  $p$ -tampering of bounded real functions). *Let  $\bar{X} = (X_1, \dots, X_n)$  be a joint distribution,  $f: \text{Supp}(\bar{X}) \mapsto [-1, +1]$  be a real function defined over  $\text{Supp}(\bar{X})$ , and  $\mu = \mathbb{E}[f(\bar{X})]$ . Then there is a  $p$ -variation  $\bar{Y}$  of  $\bar{X}$  (generated by a possibly unbounded  $p$ -tampering algorithm) such that  $\mathbb{E}[f(\bar{Y})] \geq \mu + \frac{p \cdot \text{Var}[f(\bar{X})]}{2-p \cdot (1-\mu)}$ .*

To prove Theorem C.1 we use a variant of the idea of so called ‘‘half-space sources’’ introduced in [36] and further used in [6, 19] by tailoring it to *real-output* functions. Before proving the theorem we compare its bound with what we achieved for the Boolean case.

**Comparison with the bias of Theorem 3.11.** Note that the bias obtained by the attacker of Theorem C.1 is the same as that of Theorem 3.11 when we restrict them to the Boolean case. However, these theorems are incomparable as Theorem 3.11 uses an efficient tampering (assuming  $\bar{X}$  is efficiently online samplable and  $f$  is efficiently computable) but the proof of Theorem C.1 directly defines a  $p$ -tampering variation without obtaining it through an actual efficient attack. Interestingly, the attacker of Theorem 3.11 uses a strategy (called the greedy algorithm in [1] and by us here as well) that could lead to arbitrarily small bias when applied to general (bounded) *real* functions [1]. So to prove an efficient version of Theorem C.1 one needs to use new ideas other than those of Theorem 3.11!

*Proof of Theorem C.1.* First note that the following proposition directly follows from Lemma 2.9.

**Proposition C.2.** *A distribution  $\bar{Y} = (Y_1, \dots, Y_n)$  is a  $\rho$ -tampering variation of the distribution  $\bar{X} = (X_1, \dots, X_n)$  iff  $\text{Supp}(\bar{Y}) \subseteq \text{Supp}(\bar{X})$  and for every valid prefix  $x_{\leq i} \in \text{ValPref}(\bar{X})$  we have*

$$\Pr[Y_i = x_i \mid x_{<i}] \geq (1 - \rho[x_{<i}]) \cdot \Pr[X_i = x_i \mid x_{<i}].$$

Now, for  $0 < h < 1$ , consider the distribution  $\bar{Y} = (Y_1, \dots, Y_n)$  such that  $\Pr[\bar{Y} = x]$  is defined to be  $\Pr[\bar{X} = x] \cdot \left(\frac{1+h \cdot f(x)}{1+h \cdot \mu}\right)$ . This is indeed a distribution because we have that (1)  $\Pr[\bar{Y} = x] \geq 0$  for every  $x \in \text{Supp}(\bar{Y})$  (because  $h < 1$  and  $f(x), \mu \geq -1$ ) and that (2) they sum up to one:

$$\sum_{x \in \text{Supp}(\bar{Y})} = \frac{1}{1+h \cdot \mu} \cdot \left( \sum_{x \in \text{Supp}(\bar{Y})} \Pr[\bar{X} = x] + h \cdot \sum_{x \in \text{Supp}(\bar{Y})} \Pr[\bar{X} = x] \cdot f(x) \right) = \frac{1+h \cdot \mu}{1+h \cdot \mu} = 1.$$

Now we can compute the average of  $f$  over  $\bar{Y}$  as follows:

$$\begin{aligned} \mathbb{E}[f(\bar{Y})] &= \sum_{x \in \text{Supp}(\bar{X})} \Pr[\bar{X} = x] \cdot \left(\frac{1+h \cdot f(x)}{1+h \cdot \mu}\right) \cdot f(x) \\ &= \frac{\mu + h \cdot \mathbb{E}[f(\bar{X})^2]}{1+h \cdot \mu} = \frac{\mu + h \cdot \mu^2 + h \cdot \text{Var}[f(\bar{X})]}{1+h \cdot \mu} = \mu + h \cdot \frac{\text{Var}[f(\bar{X})]}{1+h \cdot \mu}. \end{aligned}$$

All we have to do is to show that  $\bar{Y}$  is a  $\frac{2 \cdot h}{h+1}$ -tampering variation of  $\bar{X}$ , because then by setting  $p = \frac{2 \cdot h}{h+1} \in (0, 1)$  we have  $h = \frac{p}{2-p} \in (0, 1)$  and the amount of bias we get would be exactly

$$h \cdot \frac{\text{Var}[f(\bar{X})]}{1 + h \cdot \mu} = \left( \frac{p}{2-p} \right) \cdot \frac{\text{Var}[f(\bar{X})]}{1 + (p/(2-p)) \cdot \mu} = \frac{p \cdot \text{Var}[f(\bar{X})]}{2-p \cdot (1-\mu)}.$$

In the following, we focus on proving that  $\bar{Y}$  is indeed a  $\frac{2 \cdot h}{h+1}$ -tampering variation of  $\bar{X}$ . For any  $x' = (x'_1, \dots, x'_n), x'' = (x''_1, \dots, x''_n) \in \text{Supp}(\bar{X})$  we have:

$$\frac{1-h}{1+h} \cdot \frac{\Pr[\bar{X} = x']}{\Pr[\bar{X} = x'']} \leq \frac{\Pr[\bar{Y} = x']}{\Pr[\bar{Y} = x'']} \leq \frac{1+h}{1-h} \cdot \frac{\Pr[\bar{X} = x']}{\Pr[\bar{X} = x'']}.$$
 (4)

For  $i \in [n]$  and  $x = (x_1, \dots, x_n)$  define the set  $\mathcal{Z}_{x,i}$  as follows

$$\mathcal{Z}_{x,i} = \{x' = (x'_1, \dots, x'_n) \in \text{Supp}(\bar{X}) \mid x'_1 = x_1, \dots, x'_i = x_i\}$$

which is the set of all samples that are equal to  $x$  in the first  $i$  blocks. By adding Inequality 4 for every sample in  $\mathcal{Z}_{x,i-1}$  and reversing it we get

$$\frac{1-h}{1+h} \cdot \frac{\Pr[\bar{X} = x'']}{\Pr[\bar{X} \in \mathcal{Z}_{x,i-1}]} \leq \frac{\Pr[\bar{Y} = x'']}{\Pr[\bar{Y} \in \mathcal{Z}_{x,i-1}]} \leq \frac{1+h}{1-h} \cdot \frac{\Pr[\bar{X} = x'']}{\Pr[\bar{X} \in \mathcal{Z}_{x,i-1}]}.$$
 (5)

By adding Inequality 5 for every  $x'' \in \mathcal{Z}_{x,i}$  we get

$$\frac{1-h}{1+h} \cdot \frac{\Pr[\bar{X} \in \mathcal{Z}_{x,i}]}{\Pr[\bar{X} \in \mathcal{Z}_{x,i-1}]} \leq \frac{\Pr[\bar{Y} \in \mathcal{Z}_{x,i}]}{\Pr[\bar{Y} \in \mathcal{Z}_{x,i-1}]} \leq \frac{1+h}{1-h} \cdot \frac{\Pr[\bar{X} \in \mathcal{Z}_{x,i}]}{\Pr[\bar{X} \in \mathcal{Z}_{x,i-1}]}.$$

Because  $\frac{\Pr[\bar{X} \in \mathcal{Z}_{x,i}]}{\Pr[\bar{X} \in \mathcal{Z}_{x,i-1}]} = \Pr[X_i = x_i \mid x_{<i}]$ , the above is equivalent to

$$\frac{1-h}{1+h} \cdot \Pr[X_i = x_i \mid x_{<i}] \leq \Pr[Y_i = x_i \mid x_{<i}] \leq \frac{1+h}{1-h} \cdot \Pr[X_i = x_i \mid x_{<i}].$$

By Proposition C.2,  $\bar{Y}$  is a  $p$ -tampering variation of  $\bar{X}$  for  $1-p = \frac{1-h}{1+h}$ , which is equivalent to  $p = \frac{2 \cdot h}{1+h}$ .  $\square$

The existence of a (possibly inefficient)  $p$ -tampering attacker of Theorem C.1 raises the intriguing question of whether this amount of bias could be achieved through an efficient  $p$ -tampering attack.

**Limitations of Possible Bias.** Now, we turn to studying the *limitation* of  $p$ -tampering attacks in the computationally unbounded setting. We show that, perhaps surprisingly, as opposed to case of  $p$ -tampering attacks over *uniform binary* alphabet, where achieving bias  $p$  is possible for the balanced case of  $\mathbb{E}[f(\bar{X})] = 0$  even efficiently [1, 39], when it comes to blockwise  $p$ -tampering over *uniform non-binary* case, even an unbounded  $p$ -tampering attacker cannot achieve bias  $+p$  in general. We write the following proposition using binary inputs with *non-uniform* distribution, but it can then be turned into a blockwise tampering over larger blocks (with uniformly distributed bits).

**Proposition C.3.** *For a given  $n \in \mathbb{N}$  and  $\mu \in (-1, +1)$ , let  $\bar{X} = (X_1 \times \dots \times X_n)$  be a product distribution over blocks such that  $X_i$  is a biased coin where  $\Pr[X_i = +1] = (\frac{1+\mu}{2})^{1/n}$  and  $\Pr[X_i = -1] = 1 - (\frac{1+\mu}{2})^{1/n}$ . For  $x = (x_1, \dots, x_n)$  let  $f(x) = 2 \cdot (\frac{x_1+1}{2}) \cdot \dots \cdot (\frac{x_n+1}{2}) - 1$ . (Thus, both  $\bar{X}$  and  $f$  are*

parameterized by  $n$ .) Note that  $\mathbb{E}[f(\bar{X})] = \mu$ . In other words, the function is defined in a way that it is equal to 1 if all the coins are 1, and it is  $-1$  otherwise. We know that  $\Pr[f(\bar{X}) = 1] = ((\frac{1+\mu}{2})^{1/n})^n = \frac{1+\mu}{2}$  which means the expected value of the function  $f$  over  $\bar{X}$  is  $\mu$ . Now, we claim that for any  $c > (1 + \mu) \cdot \ln(\frac{2}{1+\mu})$ , there exists  $n_0 \in \mathbb{N}$  and  $p_0 \in [0, 1]$  such that for every  $p < p_0$  and  $n > n_0$  and every  $p$ -tampering variation  $\bar{Y}$  of  $\bar{X}$  we have  $\mathbb{E}[f(\bar{Y})] < \mathbb{E}[f(\bar{X})] + c \cdot p$ .

*Proof of Proposition C.3.* Since  $X_i = -1$  for any  $i$  implies  $f(\bar{X}) = -1$ , the optimal  $p$ -tampering algorithm for biasing  $f$  towards  $+1$  is to make every ‘coin’  $X_i$  to be 1 upon any chance. Let  $t = (\frac{1+\mu}{2})^{1/n}$ . Thus, the  $p$ -tampering variation  $\bar{Y} = (Y_1 \times \dots \times Y_n)$  of  $\bar{X}$  such that every  $Y_i$  is a biased coin where  $\Pr[Y_i = +1] = (1 - p) \cdot t + p$  has the highest average among all  $p$ -tampering variations of  $\bar{X}$ . It holds that

$$\mathbb{E}[f(\bar{Y})] = 2 \cdot ((1 - p) \cdot t + p)^n - 1 = 2 \cdot (1 - (1 - p)(1 - t))^n - 1 \leq 2 \cdot e^{-(1-p) \cdot (1-t) \cdot n} - 1.$$

We also have  $\lim_{n \rightarrow \infty} (1 - t) \cdot n = \lim_{n \rightarrow \infty} (1 - ((1 + \mu)/2)^{1/n}) \cdot n = -\ln((1 + \mu)/2)$ .

We are interested to know the range of the constant  $c$  when we write the amount of bias achieved by the optimal algorithm in the form  $\mu + p \cdot c$ . As we will see, for small  $p$  this constant is actually smaller than one. In particular, we have

$$\lim_{p \rightarrow 0} \lim_{n \rightarrow \infty} \frac{2 \cdot e^{-(1-p) \cdot (1-t) \cdot n} - 1 - \mu}{p} = \lim_{p \rightarrow 0} \frac{2 \cdot e^{(1-p) \cdot \ln(\frac{1+\mu}{2})} - 1 - \mu}{p} = (1 + \mu) \cdot \ln\left(\frac{2}{1 + \mu}\right).$$

Therefore, for any  $c > (1 + \mu) \cdot \ln(\frac{2}{1+\mu})$  there exist  $n_0 \in \mathbb{N}$  and  $p_0 \in [0, 1]$  such that if  $p < p_0$  and  $n > n_0$  then  $\mathbb{E}[f(\bar{Y})] < \mu + c \cdot p$ .  $\square$

Proposition C.3 shows that even  $f$  is Boolean, blockwise  $p$ -tampering attacks, cannot achieve bias  $c \cdot p$  for all  $p$  and some  $c > (1 + \mu) \cdot \ln(\frac{2}{1+\mu})$ . For the case of balanced functions where  $\mu = 0$ , it means that we cannot achieve bias  $c \cdot p$  for  $c > \ln(2) \approx 0.69$ . This leaves open the search for finding optimal constant  $c$  for which  $p$ -tampering attacks can always get at least  $c \cdot p$  bias. For case of balanced Boolean functions we already know that  $0.5 \leq c \leq \ln(2) < 0.7$ .

Also note that one can get a similar result to that of Proposition C.3 for *uniform non-binary* blocks as follows. Let  $X_i$  be distributed as  $U_m$  and  $g: \{0, 1\}^m \mapsto \{+1, -1\}$  be such that  $\mathbb{E}[g(U_m)] \approx (\frac{1+\mu}{2})^{1/n}$ . Finally, use  $g(X_i)$  instead of  $X_i$  in the definition of  $f$  in Proposition C.3.