# Doubly-Efficient Interactive Proofs for Distribution Properties

Tal Herman*  
Weizmann Institute of Science

Guy N. Rothblum†  
Apple

November 1, 2023

## Abstract

Suppose we have access to a small number of samples from an unknown distribution, and would like learn facts about the distribution. An untrusted data server claims to have studied the distribution and makes assertions about its properties. Can the untrusted data server *prove* that its assertions are approximately correct? Can a short efficiently verifiable proof be generated in polynomial time?

We study *doubly-efficient* interactive proof systems that can be used to verify properties of an unknown distribution over a domain $[N]$. On top of efficient verification, our focus is on proofs that the honest prover can generate in polynomial time. More generally, the complexity of generating the proof should be as close as possible to the complexity of simply running a standalone analysis to determine whether the distribution has the property.

Our main result is a new 2-message doubly-efficient interactive proof protocol for verifying any label-invariant distribution property (any property that is invariant to re-labeling of the elements in the domain of the distribution). The sample complexity, communication complexity and verifier runtime are all $\widetilde{O}(\sqrt{N})$. The proof can be generated in *quasi-linear* $\widetilde{O}(N)$ time and sample complexities (the runtimes of the verifier and the honest prover hold under a mild assumption about the property's computational complexity). This improves on prior work, where constructing the proof required super-polynomial time [Herman and Rothblum, STOC 2022]. Our new proof system is directly applicable to proving (and verifying) several natural and widely-studied properties, such as a distribution's support size, its Shannon entropy, and its distance from the uniform distribution. For these (and many other) properties, the runtime and sample complexities for generating the proof are within $\mathsf{polylog}(N)$ factors of the complexities for simply determining whether the property holds.

# Contents

# 1 Introduction

Given sample access to an unknown discrete distribution over a large domain $[N]$, what can we learn about the distribution's properties? How many samples are required, and what is the computational complexity of learning? These are foundational questions in statistics and in computer science.

An emerging line of works asks a new question: what is the complexity of *verifying* facts about the distribution? Suppose that an *untrusted* prover, who knows the distribution, claims that the distribution has some property, e.g. that the distribution's entropy is $k$, or that its support size is $M$. Can the prover provide a proof of approximate correctness for such claims? We are interested in proofs that can be verified using fewer samples and computational resources than it would take to approximate these quantities on our own. More generally: which distribution properties can be verified efficiently?

The computational complexity and sample complexity *of generating the proof* are also of central importance. This is true both from a foundational perspective, and with the goal of deploying proof systems in the real world, where generating the proof needs to be as efficient as possible, and certainly needs to be computationally feasible. In a *doubly-efficient proof system* the proof can be generated in polynomial time. More generally, the complexity of proving should be as close as possible to the complexity of simply performing the task. In this work we study doubly-efficient proof systems for verifying properties of distributions, our main question is:

*Can an untrusted prover generate,* in polynomial time, *a proof that convinces a verifier that an unknown distribution has some property? How efficient can proving and verifying be?*

We focus on verifying distribution properties via an interactive proof system [GMR85], where a probabilistic verifier has sampling access to the distribution and communicates with an untrusted prover. This continues a study of proof systems for distribution properties initiated by Chiesa and Gur [CG18]. Drawing inspiration from the property testing literature [GGR98, RS96], the prover's claim is that the distribution has (or is close to having) a property. If the prover's claim is approximately correct, the verifier accepts with high probability. If the claim is *far* from correct, i.e. the distribution is far from the property, then no matter what strategy a cheating prover might follow, the verifier rejects with high probability. Recently, Herman and Rothblum [HR22] showed that, for the rich class of "label invariant" distribution properties (see below), *approximate verification can be very efficient*: any such property has a 2-message interactive proof system, where verification requires only $\tilde{O}(\sqrt{N})$ samples, communication, and (under a mild additional assumption) verification time. These results apply to several important and widely-studied tasks, such as estimating the Shannon entropy, the support size, and the distance from the uniform distribution, which all boil down to verifying label-invariant distribution properties. On the other hand, for all these properties, performing the analysis (without help from an untrusted prover) requires $\Omega(N/\log N)$ samples and running time [VV10]. Thus, for these (and other) properties, verification can be quasi-quadratically more efficient than performing the analysis.

While this result showed the *existence* of efficiently-verifiable proofs for a rich class of data analyses, the computational complexity of generating the proof was prohibitive: the honest prover runtime was *super-polynomial* in $N$. This made their protocol infeasible for many scenarios. Moreover, for the important and widely-studied properties mentioned above, there was a huge gap between the complexity of *computing* (or deciding) the problem (which is quasi-linear in $N$) and the complexity of *generating the proof* (which was super-polynomial).

## 1.1 This Work: Doubly-Efficient Proofs

Our main result is a new *doubly-efficient* proof system for approximately verifying any label-invariant property of distributions over the domain $[N]$, where the honest prover's runtime is *quasi-linear in N*. For the support size, the Shannon entropy and the distance from the uniform distribution, the proof can be generated with only $\mathsf{polylog}(N)$ overhead over the computational complexity and the sample complexity of simply deciding the problem (without generating a proof).

We proceed to detail this result: a *distribution property* is a set of distributions (similarly to the way a *language* is a set of strings), parameterized by the size of the domain $N$. *Label-invariant* properties (sometimes referred to as *symmetric* properties) are a natural class of distribution properties, where changing the "labels" of elements in the support of a distribution does not change membership in the property.[1] Many natural and widely-studied properties are label-invariant: e.g. being uniform over $[N]$, having entropy $k$, or having support size $M$. Another example is being uniform on a set of size $S$ [BC17]. On the other hand, having support on the odd elements in $[N]$ is (one example of) a property that is *not* label-invariant. We measure the distance of a distribution $D$ from a property $\mathcal{P}$ by $D$'s total variation distance to the closest distribution in $\mathcal{P}$.

**Theorem 1.1** (Main result: doubly-efficient IPs for label-invariant properties). *For every label-invariant property $\mathcal{P}$ with a doubly-efficient approximate decision procedure (see below), there exists an interactive proof system as follows. The prover and the verifier both get as input an integer $N$ and proximity parameters $\varepsilon_c, \varepsilon_f \in [0, 1]$ where $\varepsilon_c < \varepsilon_f$, as well as sampling access to an unknown distribution $D$ over support $[N]$, and the following properties hold:*

- *Completeness: if $D$ is $\varepsilon_c$-close to the property (i.e. $D$ is at statistical distance at most $\varepsilon_c$ from a distribution that has the property) and the prover follows the protocol, then the verifier accepts w.h.p.*

- *Soundness: if $D$ is $\varepsilon_f$-far from the property (its statistical distance from every distribution in the property is at least $\varepsilon_f$), then no matter how the prover cheats, the verifier rejects w.h.p.*

- *Doubly-efficient prover: Taking $\rho = \varepsilon_f - \varepsilon_c$, the honest prover's runtime and sample complexity are $\widetilde{O}(N) \cdot \mathsf{poly}(1/\rho)$.*

- *Efficient verification: the protocol consists of 2 messages, the communication complexity and the verifier's sample complexity and runtime are all $\widetilde{O}(\sqrt{N}) \cdot \mathsf{poly}(1/\rho)$.*

We emphasize that the completeness requirement is *tolerant* [PRR06]: the verifier should accept even if the distribution is not in the property, so long as it is *close to the property*. The complexity is polynomial in the gap ($\varepsilon_f - \varepsilon_c$) between the distances. Tolerant verification can be used to approximately verify the distribution's distance to the property: if the prover claims the distance is $\delta$, we can verify this (up to distance $\rho$) by setting $\varepsilon_c = \delta$ and $\varepsilon_f = \delta + \rho$ in our proof system.

**Doubly efficient approximate decision condition.** We need the property to satisfy a mild *approximate decision condition*. This assumes the existence of two procedures as follows. The first procedure, given a histogram of the distribution's probabilities, accepts if the distribution is in

---

[1]More formally, for a distribution $D$ over the domain $[N]$, and a permutation $\pi : [N] \rightarrow [N]$, we let $\pi(D)$ be the distribution obtained by sampling from $D$ and applying the permutation $\pi$ to the outcome. A property $\mathcal{P}$ is *label-invariant* if for every distribution $D \in \mathcal{P}$, and every permutation $\pi$ over $D$'s domain, $\pi(D) \in \mathcal{P}$.

the property and rejects if the distribution is $\sigma$-far from the property. The histogram is defined as follows: we bucket the elements in the domain according to their (approximate) probabilities, and the histogram specifies the number of elements in each bucket. In more detail, we round each element's probability down to the nearest value $e^{\ell\tau}/N$ for an integer $\ell$, where $\tau$ is an approximation parameter that is polynomial in $\sigma$. We refer to this as the distribution's $\tau$-*approximate histogram*, and note that it gives sufficient information for approximating the distribution's distance from a label-invariant property. We can ignore the elements whose probabilities are very small, so the $\tau$-approximate histogram can be represented using $O(\log^2 N/\tau)$ bits. The approximate decision procedure should run in time that is polynomial in this representation, i.e. in $\mathsf{poly}(\log N, 1/\sigma)$ time (see Definition 6.5), and is used by the verifier in our interactive proof.

We also need a second procedure that, given sample access to a distribution $D$ that is $\varepsilon_c$-close to the property, outputs the histogram of a distribution $D'$ in the property that is (approximately) $\varepsilon_c$-close to $D$. This second procedure is used by the honest prover in our protocol, and we require that it runs in quasi-linear time (more generally, the honest prover runtime grows with the runtime of this procedure). A property that has both procedures has a *doubly-efficient approximate decision procedure*. We view this as a mild condition, and note that it is satisfied by natural properties such as the support size, the distance from $U_N$ and the Shannon entropy.

**Verified histogram and applications.**    The protocol of Theorem 1.1 is achieved by constructing a sub-protocol that lets the verifier learn a *verified (approximate) histogram* of the distribution $D$. The probability that the verifier accepts and the histogram is $\sigma$-far from accurate for $D$ is small, where the complexities are all polynomial in $1/\sigma$. In turn, the verified histogram can be used to obtain the result of Theorem 1.1, estimating the distance from a label-invariant property, or to obtain protocols for quantities of interest that can be estimated from the histogram.

In particular, we get *doubly-efficient* protocols for the quantities discussed above:

- **Distance to $U_N$:** given claimed distance $\delta$, if $|\Delta(D, U_N) - \delta| > \rho$, the verifier rejects w.h.p.

- **Support size:** given a promise that each element in $D$'s support has probability at least $1/N$ and claimed support size $M$, if $||\mathrm{Supp}(D)| - M| > \rho \cdot N$, the verifier rejects w.h.p. The promise of a lower bound on the probabilities of elements in D's support is standard in the study of estimating the support size.

- **Shannon entropy:** given claimed entropy $k$, if $|\mathsf{H}(D) - k| > \rho$, the verifier rejects w.h.p.

For all these problems, if the prover's claim is (approximately) correct, then the verifier accepts w.h.p. The proof system for distance from $U_N$ follows immediately from obtaining a verified histogram. The proof systems for the entropy and support size also follow by translating the statistical distance between the distribution $D$ and the claimed histogram into a bound on the difference between $D$'s entropy or support size and the value implied by the histogram. The complexities of all protocols are polynomial in $(1/\rho)$ (for the Shannon Entropy we need the verified histogram to be $(\rho/\log N)$-accurate w.r.t to the distribution $D$).

We remark that our protocol actually gives a stronger guarantee: the verifier can obtain a collection of samples, drawn i.i.d. from $D$, and alleged probabilities for each of these samples, where the probabilities are guaranteed to be approximately correct (see Remark 2.2). The verified tagged samples can be used to derive an approximate histogram (deriving the size of each bucket $\ell$ from the fraction of samples tagged as belonging to bucket $\ell$ divided by the probability of elements in that bucket), but they may also have further applications.

**On our protocols' complexities.** Several remarks about the protocol's complexity are in order. The sample complexity is nearly-optimal (for any interactive proof, regardless of its communication or round complexities): Chiesa and Gur [CG18] (extended in [HR22]) showed an $\Omega(\sqrt{N})$ sample-complexity lower bound for the promise problem where in the YES case the distribution equals $U_N$ and in the NO case the distribution is uniform over a set of size $(N/2)$. Thus, this lower bound applies also to verifying the distance from uniform, the Shannon entropy, and the support size.

Our protocols all use secret coins. For clarity of exposition, our protocols are presented as if the honest prover has *perfect* knowledge of the distribution, but this idealized honest prover can implemented by a quasi-linear time honest prover that learns a sufficiently-accurate (multiplicative) approximation to the distribution. As noted above, the sample and runtime complexities of standalone distribution testing (without generating a proof) for the distance form uniform, entropy, and support size properties are $\Theta(N/\log(N))$ [VV10]. Thus, the prover complexity is optimal for these quantities, and is within $\mathsf{polylog}(N)$ factors of the complexity of deciding the problem.

**Comparison to known results.** Our result is most directly related to the interactive proof of [HR22]. Our main contribution is achieving a quasi-linear in $N$ runtime for the honest prover, whereas in the prior work the honest prover runtime was $N^{\log(N) \cdot \mathsf{poly}(1/\rho)}$ (i.e., super-polynomial).

We also compare to two other known methods for verifying general distribution properties, which are both doubly-efficient. First, Chiesa and Gur [CG18] showed it is possible to verify using small sample complexity but large communication and verifier runtime. In their protocol, the prover sends a complete description of a distribution $\widetilde{D}$. The verifier checks that $\widetilde{D}$ is close to the property, and then runs a distribution tester to verify that the alleged distribution $\widetilde{D}$ is $\varepsilon$-close to the actual distribution $D$. The verification can be performed using $O(\sqrt{N}/\varepsilon^2)$ samples [BFF+01, VV14, Gol20]. Moreover, the protocol is non-interactive, using only a single message. However, the verification time and the communication are *quasi-linear in $N$*. It is also possible to verify with zero communication by having the verifier ignore the prover and simply learn (an approximation to) the entire distribution $D$ on its own (see Theorem 3.15). This requires no communication, but the sample complexity and verification time are linear in $N$.

In contrast to the above solutions, our focus is on verification that is simultaneously efficient in terms of the verifier's running time, of the communication complexity, and of the sample complexity. In our protocols, all of these complexity measures are bounded by $\widetilde{O}(\sqrt{N}) \cdot \mathsf{poly}(1/\sigma)$.

## 1.2 Further Related Work

We study the verification of distribution properties via interactive proofs. Interactive proof systems were introduced by Goldwasser, Micali and Rackoff [GMR85] in the context of proving computational statements about an input that is fully known to the prover and the verifier. In our work, the distribution can be thought of as the input, but it is not fully known to the verifier. We aim for verification without examining the distribution in its entirety, using minimal resources (samples, communication, runtime, etc.). Our work builds on a line of work that studies the power of sublinear time verifiers, who cannot read the entire input [EKR04, RVW13, GR18], on verifying properties of distributions using a small number of samples [CG18, HR22], and on verifying the result of machine learning algorithms using a small number of labeled examples [GRSY21]. In particular, Chiesa and Gur [CG18] introduced and studied interactive proofs for distribution verification and showed upper and lower bounds. Our work is most closely related to (and builds on) the protocol of [HR22] for label-invariant properties, where proof generation required super-polynomial time.

**Doubly-efficient proof systems.** Our work focuses on doubly-efficient proof systems, where the honest prover can generate the proof in polynomial time. This is motivated by the goal of building interactive proof systems that can be used in the real world (where all parties, including the honest prover, need to run in polynomial time). It is also very important from a foundational perspective, where achieving polynomial runtime for the honest prover is of central importance. This was already true in the genesis of the field: the prover in a zero-knowledge proof for an NP language is required to run in polynomial time given a witness to the input's membership in the language [GMR85, GMW91]. It was also an early focus in works on PCPs [BFLS91], on computationally sound CS proofs [Mic94], and in the line of work on doubly-efficient interactive proof systems for delegating computation [GKR15].

## 2 Technical Overview

### 2.1 The Protocol Behind Theorem 1.1

We describe the protocol behind Theorem 1.1 in broad strokes , and review several of our technical ideas and contributions.

Membership in a label-invariant property can be decided based on the probability histogram of the distribution, i.e. for every $p \in [0,1]$ *how many* elements $x \in [N]$ satisfy $D(x) = p$. Through our protocol, the verifier obtains an approximation of this object, namely, the *(approximate) bucket-histogram* of the distribution $D$, that has description of size $\mathsf{polylog}(N)$ bits (compared to the potentially $\Omega(N)$ size of the probability histogram). The bucket histogram allows the verifier to approximate the distance of $D$ from any label-invariant property.

For an accuracy parameter $\tau < 0.01$[2], the $\tau$-*approximate bucket histogram* partitions the interval $[0,1]$ into $O(\log N/\tau)$ *buckets*, and counts how much *probability mass* of the distribution $D$ falls in each bucket. More concretely, we define the $\ell$'th bucket of $D$ to be:

$$B_\ell^D = \left\{ x \in [N] : D(x) \in \left[ \frac{e^{\ell\tau}}{N}, \frac{e^{(\ell+1)\tau}}{N} \right) \right\}$$

And denote its mass $p_\ell = D(B_\ell^D)$. We consider all elements $x \in [N]$ with probability, $D(x) \leq \frac{\tau^2}{N}$ to be in one bucket, with corresponding index $L$. We call the collection $\{p_\ell\}_{\ell \in \mathbb{Z}: \frac{e^{\ell\tau}}{N} \geq \frac{\tau^2}{N}}$ the $\tau$-approximate bucket histogram of $D$. (We omit the subscript from now on for ease of notation.)

**Assuming $D$ contains no *heavy* elements.** Similar to [HR22], we show a protocol to obtain the bucket histogram of the distribution $D$ assuming that it has no *heavy* elements, that is, for all $x$, $D(x) = O\left(\frac{1}{\sqrt{N}}\right)$. In a nutshell, by taking $\widetilde{O}(\sqrt{N})$ samples from $D$, it is possible to approximate well enough the probability of every element $x$ with probability $D(x) = \Omega\left(\frac{1}{\sqrt{N}}\right)$, and so, in order to obtain the full bucket histogram of the distribution, all that is left is to compute the histogram on the *lighter* part of the domain, which is what our protocol achieves.

---

[2]In the context of tolerant verification of label-invariant distribution properties, i.e. looking to accept distributions $\varepsilon_c$ close to the property and rejecting distributions $\varepsilon_f$ far from it (in total variation distance), we take $\tau \approx (\varepsilon_f - \varepsilon_c)^2$

### 2.1.1 Approximating the bucket histogram of a distribution with no heavy elements

**Communication Phase.** The verifier performs the following sampling process $2s$ times for $s = \widetilde{O}\left(\sqrt{N}\right)\mathsf{poly}(\tau^{-1})$: it flips a fair coin and obtains $b \in \{0,1\}$, if $b = 0$, it draws $z$ by $D$, and if $b = 1$ it draws $z$ by $U_{[N]}$ (the uniform distribution over the entire domain).

The verifier thus obtains the bits $b \in \{0,1\}^{2s}$, and the respective sample, $(z_i)_{i\in[2s]}$, which is composed of two intertwined samples, one from $D$ and one form $U_{[N]}$, denoted by $S_D = \{i \in [2s] : b_i = 0\}$ and $S_U = \{i \in [2s] : b_i = 1\}$, each of size roughly $s$.

The verifier sends the sample $(z_i)_{i\in[2s]}$ to the prover, who replies with the *tag* of each sample: the alleged bucket index to which the element $z_i$ belongs. That is, the prover sends $(\mathrm{tag}(z_i))_{i\in[2s]}$, such that, allegedly, $D(z_i) \approx \frac{e^{\mathrm{tag}(z_i)\tau}}{N}$. If $z_i \notin \mathrm{Supp}(D)$, the prover sends $\perp$. The verifier then computes the alleged empirical mass of bucket $j$ according to the sample $S_D$, $v_j = \frac{|\{i\in S_D : \mathrm{tag}(z_i)=j\}|}{|S_D|}$.

Observe that if the prover is honest, then $v_j \approx p_j$ for every $j$, and the verifier obtained a good approximation of the approximate bucket histogram of $D$! The rest of the protocol involves verifying that indeed the prover didn't lie. In order to do so, the verifier performs two tests, which are carried out without any further interaction with the prover.

**Test 1: bucket size verification.** First, The verifier checks that no element sampled by $D$ was tagged $\perp$. Then, they check that the alleged empirical mass of the $j$'th bucket as observed in the part of the sample drawn by $D$ matches the expected *size* of the bucket, reflected through the samples drawn by $U_{[N]}$.

Concretely, we expect the empirical mass of the $j$'th bucket according to the samples drawn from $U_{[N]}$, to be roughly $\frac{|B_j^D|}{N}$. The verifier doesn't know $\left|B_j^D\right|$, but it knows $v_j$, which should be close to the mass of the $j$'th bucket according to $D$. And so, the verifier computes $\frac{v_j}{e^{j\tau}/N}$ and uses it as an approximation of the alleged quantity $\left|B_j^D\right|$. For every $j$, the verifier counts how many samples in $S_U$ were tagged $j$, and computes from that the alleged empirical mass of the $j$'th bucket *as observed from the samples drawn from* $U_{[N]}$. They then check that this quantity is roughly $v_j e^{-j\tau}$.

Note that if the prover is honest $v_j e^{-j\tau} \approx p_j e^{-j\tau} = \frac{Np_j e^{-j\tau}}{N} \approx \frac{|B_j^D|}{N}$.

**Test 2: Collisions Matching Test.** For every alleged bucket $j$, let the set of samples in $S_D$ tagged $j$ be $S_D^j$. Note that by definition $\left|S_D^j\right| = sv_j$. The verifier expects that the true mass according to $D$ of the set $\{z_i\}_{i\in S_D^j}$ is roughly $\left|S_D^j\right| \cdot \frac{e^{j\tau}}{N}$.[3] And so, the verifier draws a fresh sample of size $s$ by $D$, and computes the empirical mass of the set $\{z_i\}_{i\in S_D^j}$ according to the new sample. If this mass is not roughly $\left|S_D^j\right| \cdot \frac{e^{j\tau}}{N}$, for all $j$, the verifier rejects.

If neither test failed, the verifier accepts.

**Completeness.** If the prover is honest, both tests pass, and we get that $v_j \approx p_j$ for all $j$. The honest prover strategy only requires them to know the probability of the elements sent by the

---

[3]Assume for sake of simplicity that $S_D^j$ contains only unique elements. In actuality, through choice of $s$ and the assumption that $D$'s support does not contain *heavy elements*, the number of collisions inside $S_D$ is small with respect to the expected error.

verifier, which can be approximated with sufficient accuracy using $\widetilde{O}(N)\mathsf{poly}(\tau^{-1})$ samples and runtime (for a detailed discussion of the complexity of the honest prover and the level of accuracy it requires see Remark 4.14).

**Soundness.** We show that no matter what strategy a prover might employ, if the prover *significantly* miss-tags the samples, then one of the tests will fail with high probability.

We characterize dishonest prover behavior by considering, for every two buckets $\ell$ and $j$, the variable $x_{\ell,j} \in [0,1]$, which is the fraction of samples that were sampled from $D$ and *truly* landed in the bucket $B_\ell^D$, but were tagged $j$. Moreover, we consider the number of samples from $D$ that truly landed in $B_\ell^D$ is $|S_D|\, p_\ell$ (recall that $p_\ell = D(B_\ell^D)$, so this is indeed the expected number of samples from bucket $\ell$ drawn in $S_D$). Thus, the number of samples in $S_D$ that landed in bucket $\ell$, but were claimed to be in bucket $j$, is $|S_D|\, p_\ell x_{\ell,j}$.

**Analyzing Test 1.** We show that even though the variables $\{x_{\ell,j}\}$ reflect how the prover lies on the samples drawn from $D$, they also capture the way the prover lies on $S_U$. That is, $x_{\ell,j}$ is also close to the fraction of samples drawn by $U_{[N]}$ that landed in $B_\ell^D$, but were claimed to belong to bucket $j$. This is due to the fact that the bits $b \in \{0,1\}^{2s}$ are hidden from the prover. Upon receiving a sample $z$, all the prover can know is to which true bucket $\ell$ the element $z$ belongs. Looking at the set of all samples in $(z_i)_{i\in[2s]}$ that landed in $B_\ell^D$, the prover knows that roughly $s \cdot p_\ell$ were sampled from $D$, and the rest, roughly $s \cdot \frac{|B_\ell^D|}{N} = s \cdot p_\ell e^{-\ell\tau}$ were sampled from $U_{[N]}$, however, the choice of which samples were drawn by $D$ and which from $U_{[N]}$ is unknown to the prover.

Therefore, a prover that wishes to miss-tag $x_{\ell,j}$ fraction of the samples in $S_D$ that fell in $B_\ell^D$ and tag them $j$, will also in the process similarly miss-tag roughly $x_{\ell,j}$ fraction of the samples drawn from $U_{[N]}$ that landed in the $\ell$'th bucket.

Focusing our attention on the samples that were tagged $j$, and taking $\frac{|B_\ell^D|}{N}$ to be (a good approximation of) the true fraction of samples that were sampled by $U_{[N]}$ and landed in $B_\ell^D$, we get that the fraction of samples drawn by $U_{[N]}$ and tagged $j$ is $\sum_\ell \frac{|B_\ell^D|}{N} \cdot x_{\ell,j} = \sum_\ell p_\ell e^{-\ell\tau} x_{\ell,j}$, where the last equality is achieved by plugging $|B_\ell^D| = N p_\ell e^{-\ell\tau}$. Test 1 checks that for every alleged bucket $j$, the fraction of samples drawn by $U_{[N]}$ and tagged $j$ equals $v_j e^{-j\tau}$. This amount to requiring that the prover's miss-tags satisfy the following equation:

$$e^{-j\tau} \approx \sum_\ell \frac{p_\ell x_{\ell,j}}{v_j} \cdot e^{-\ell\tau} \tag{1}$$

**Remark 2.1.** *Note that heavy buckets, i.e. buckets $\ell$ for which $\frac{e^{\ell\tau}}{N} > \frac{1}{\sqrt{N}}$, might have considerable mass, yet very few elements. Upon receiving some sample $z$ from a heavy bucket, the probability that they were sampled by $U_{[N]}$, and not through $D$ is roughly $\frac{1/N}{e^{\ell\tau}/N} < \frac{1}{\sqrt{N}}$, and so, the prover can tag roughly $O(\sqrt{N})$ such samples, and still not miss-tag any sample that was drawn by $U_{[N]}$ and fell in $B_\ell^D$. Limiting our scope to distributions with no heavy elements, as well as taking sample complexity $s = \widetilde{O}\left(\sqrt{N}\right)$ assures us that for any bucket with significant mass, we can establish a relation between the cheating patterns across $S_D$ and $S_U$.*

**Analyzing Test 2.** For every alleged bucket $j$, the true mass of *the set of samples tagged $j$* is $\sum_{i\in S_D^j} D(z_i)$. And so, the expected empirical mass of $S_D^j$ in the fresh sample should be roughly

7

$\sum_{i \in S_D^j} D(z_i)$. Since this sum is composed of elements from potentially many other buckets, we rewrite the sum to reflect this. There are roughly $sp_\ell x_{\ell,j}$ samples that fell in bucket $\ell$ (i.e. each with true probability $\frac{e^{\ell\tau}}{N}$) and were tagged $j$, and so $\sum_{i \in S_D^j} D(z_i) = \sum_\ell sp_\ell x_{\ell,j} \frac{e^{\ell\tau}}{N}$.

Since Test 2 checks that the empirical mass of $\{z_i\}_{i \in S_D^j}$ in $S_U$ is roughly $\left| S_D^j \right| \cdot \frac{e^{j\tau}}{N} = sv_j \frac{e^{j\tau}}{N}$, we get that Test 2 essentially verifies that:

$$e^{j\tau} \approx \sum_\ell \frac{p_\ell x_{\ell,j}}{v_j} \cdot e^{\ell\tau} \tag{2}$$

**Combining tests in "sterile" setting.** Assume that Approximate Equations (1) and (2) hold with exact equality. For every $j$, consider the distribution $P_j$ that assigns every bucket index $\ell$ the probability $\frac{p_\ell x_{\ell,j}}{v_j}$. In this case, rewrite the equations as follows:

$$\mathbb{E}_{\ell \sim P_j} \left[ e^{\ell\tau} \right] = e^{j\tau}$$

$$\mathbb{E}_{\ell \sim P_j} \left[ e^{-\ell\tau} \right] = e^{-j\tau}$$

By Jensen's Inequality: $e^{j\tau} = \mathbb{E}_\ell \left[ e^{\ell\tau} \right] \geq e^{\tau\mathbb{E}[\ell]}$, and $e^{-j\tau} = \mathbb{E}_\ell \left[ e^{-\ell\tau} \right] \geq e^{-\tau\mathbb{E}[\ell]}$, or equivalently, $e^{j\tau} \leq e^{\tau\mathbb{E}[\ell]}$. We conclude that $e^{\tau\mathbb{E}[\ell]} = e^{j\tau}$, from which we get that

$$\mathbb{E}_\ell \left[ e^{\ell\tau} \right] = e^{\tau\mathbb{E}[\ell]}$$

Meaning that Jensen's Inequality holds with equality for every $j$. This can only happen if random variable $e^{\ell\tau}$ is a constant in $P_j$, i.e. that $x_{\ell,j} = \mathbb{1}_{\ell=j}$. Setting $x_{\ell,j} = \mathbb{1}_{\ell=j}$ for all $\ell$ and $j$ is exactly the honest prover strategy, where every sample is tagged correctly. In other words, the only prover strategy that satisfies both tests for all $j$ is the honest prover behavior, and any other strategy will be rejected.

Of course, we don't expect Equations (1) and (2) to hold in exact equality upon running the actual protocol. In the next section we show that assuming these equations are *close* to be correct is enough to argue that any prover response that's *far enough* from the honest strategy will be rejected with high probability.

**Remark 2.2** (Verified tagged sample). *Our protocol guarantees even more than approximate correctness of the histogram. The set of taggings on the samples drawn from $D$ is guaranteed to be "close" to the true probabilities of those samples (otherwise the verifier rejects w.h.p.). This is a potentially more powerful guarantee that may lead to further applications, e.g. in testing for properties of pairs of distributions. See Theorem 6.1 for the formal guarantee.*

### 2.1.2 Slack analysis: single bucket case

We show how to account for slack by focusing on a simplified case, where the prover claims that the entire distribution $D$ is supported on a single bucket, $j_0$, while $D$ is actually $\sigma$-far from it in total variation distance (see Definition 3.1).

In other words, if we denote by $S$ the alleged support size of $B_{j_0}^D$, we'd want to accept when $D$ is uniform over $S$ elements, and reject if it is $\sigma$-far from that. This problem was investigated in the non-interactive setting in [BC17]. Note that $S \in \left( Ne^{-(j_0+1)\tau}, Ne^{-j_0\tau} \right]$.

**Test 1: Revisited.** The verifier checks that for every $i \in S_D$, $\text{tag}(z_i) = j_0$, i.e. that no sample that was drawn from $D$ was claimed to be outside $\text{Supp}(D)$. Next, set $w_{j_0}$ to be the fraction of samples tagged $j_0$ in $S_U$. Recall that assuming that the $D$ is not supported on *heavy* buckets means that $S = \Omega(\sqrt{N})$, and through our choice of $s = \widetilde{O}(\sqrt{N})\text{poly}(\tau^{-1})$, the verifier expects the empirical mass of $\text{Supp}(D)$ in $S_U$ to be in $(1 \pm O(\tau))\frac{S}{N}$ with high probability. However, if $\text{Supp}(D) \notin (1 \pm O(\tau)) S$, then the prover has to miss-tag elements in order to pass this test.

Consider the case where $D$ is $\sigma$-far from being uniform over $S$ elements *and also* $|\text{Supp}(D)| = (1 + \Omega(\tau)) S$. Then, in order to pass the test with high probability, the prover must claim that samples $z$ that fell inside the support of $D$, do not belong to the support. However, by doing so, the prover risks miss-tagging a sample drawn from $D$, and failing the first part of the test. Still, the prover can miss-classify and pass the test with high probability, by choosing to miss-classify elements with low probability according to $D$, which are much more likely to have been drawn from $U_{[N]}$.

Assume therefore that *$D$'s support doesn't contain any tiny probability element.* I.e. for all $x \in \text{Supp}(D)$, $D(x) \geq \frac{\tau}{N}$ (we later discuss what happens when there are tiny elements). In this case, when the prover receives a sample $z \in \text{Supp}(D)$, the probability that the sample was drawn from $D$ is $\frac{D(z)}{D(z)+1/N} \gtrsim \tau$, and so, when the prover miss-classifies $\Omega(1/\tau)$ samples, it is likely that at least one of the samples was drawn from $D$, and the verifier will reject. We conclude that the prover can only miss-label $O(\tau^{-1})$ samples as being outside the support and still pass Test 1 with high probability. Therefore, if Test 1 passed, with high probability:

$$|\text{Supp}(D)| \leq (1 + O(\tau)) S + \frac{O(\tau^{-1})}{s} \cdot N \tag{3}$$

**Test 2: Revisited.** In the one alleged bucket case, Test 2 amounts to drawing a fresh sample of size $s$ by $D$, computing the empirical mass of all the elements in $\{z_i\}_{i \in S_D}$, and then comparing it to their the expected mass, which should be $|S_D|\frac{e^{j_0\tau}}{N} \approx s \cdot \frac{e^{j_0\tau}}{N}$. For any $D$, the expected empirical mass of $S_D$ in the fresh sample will be $s \cdot \sum_{i \in S_D} D(z_i)$. Taking the expectation also over the choice of $S_D$, we get that with high probability, this quantity is up to a multiplicative factor of $\tau$ close to $s \cdot \|D\|_2^2$. Therefore, if Test 2 passed, then with high probability $\|D\|_2^2 \in (1 + \tau) \frac{1}{S}$.

If $D$ is also $\sigma$-far in statistical distance from being uniform over $S$ elements, Herman and Rothblum [HR22] prove the following lemma:

**Lemma 2.3** (Support Size Gap Lemma [HR22]). *For every discrete distribution $D$, integer $S \in \mathbb{N}$, and parameters $\sigma, \tau \in [0,1]$. If $D$ satisfies:*

- $\|D\|_2^2 \in \left[\frac{1-\tau}{S}, \frac{1+\tau}{S}\right]$,

- *$D$ is at statistical distance at least $\sigma$ from every distribution that is uniform over $S$ elements,*

*then:*

$$|Supp(D)| \geq S \left(1 + O(\sigma^2) - O(\tau)\right) \tag{4}$$

If $D$ is $\sigma$-far from being uniform over $S$ elements, and have passed Test 1 and Test 2, then, with high probability, both Inequalities (3) and (4) hold. We show that there exists a choice of $\tau = O(\sigma^2)$ and $s = \widetilde{O}\left(\sqrt{N}\right)\text{poly}(\tau^{-1})$, as well as constants $0 < c_1 < c_2$ such that if $D$ is $\sigma$-far from

9

uniform, only distributions satisfying $\mathrm{Supp}(D) \leq (1 + c_1) S$ can pass Test 1 with high probability, while only distributions that satisfy $\mathrm{Supp}(D) \geq (1 + c_2) S$ can pass Test 2 with high probability; and so at least one of the tests must fail.

**Remark 2.4** (Soundness when $D$ contains *tiny* elements.). *Assume that $C = \left\{ x \in Supp(D) : D(x) < \frac{\tau}{N} \right\}$ isn't empty. Note that in this case, it might that the prover can miss-tag many samples in the support of $D$ as being outside the support, without landing on a sample drawn by $D$, and failing Test 1. This is since some samples might be considerably less likely to have been sampled by $D$ than by $U_{[N]}$. And so, instead of thinking of Test 1 and Test 2 as producing two conflicting claims about $|Supp(D)|$, we think of them doing so about $|Supp(D) \setminus C|$.*

## 2.2 Comparison with the [HR22] Protocol

Herman and Rothblum [HR22] give a protocol for obtaining the approximate histogram of a samplable distribution $D$. However, the construction suffered from super-polynomial prover runtime. In order to discuss the differences between the constructions, we first present the approach of [HR22] to the verification problem of the previous section. Namely, a protocol for accepting distributions uniform over a set of size $S$, and rejecting distributions $\sigma$-far from any distribution uniform over $S$ elements.

The verifier computes a a $\tau$-approximation to $D$'s collision probability, $\|D\|_2^2$, in the same way presented in the previous section. Then, through a variant of Lemma 2.3, they show that if Test 2 passed, then the Shannon entropy of any $D$ which is $\sigma$-far from uniform over a set of size $S$, must satisfy $H(D) \geq \log S + \frac{\sigma^2}{32} - \tau$.

Thus, in this case, $D$'s true entropy is significantly *higher* than what is claimed by the cheating prover (the prover claims the entropy is $\log(S)$, and for $\tau = O(\sigma)$, there is a $\Theta(\sigma^2)$ entropy gap). To detect this false claim, in the [HR22] protocol, the verifier asks the prover to execute an *entropy upper bound protocol* from the statistical zero knowledge literature [SV03, Vad99]. This protocol requires super-polynomial running time for the honest prover: see Section 2.2 for further elaboration and a comparison to our approach. We want a doubly-efficient proof system, so using the entropy upper bound protocol is a non-starter: we need a new approach.

we elaborate briefly on the entropy upper bound protocol, which is the source for the honest prover's super-polynomial running time in their work. The first step in the entropy upper bound protocol is amplifying the Shannon Entropy gap into a *min-entropy* gap by repetition (sometimes referred to as a "direct product"): taking many samples from the distribution $D$. This canonical idea goes back to the work of [HILL99]. Taking $m = \Theta(\log(N)/\sigma^2)$ samples from $D$ gives a new distribution $D^{\otimes m}$ with *min-entropy* close to $(m \cdot \mathsf{H}(D))$. The min-entropy of $D^{\otimes m}$'s can be upper bounded via a standard protocol that uses a strong seeded randomness extractor (see e.g. Vadhan [Vad12]). In the NO case, where $D$'s Shannon entropy was at least $(\log(S) + \Theta(\sigma^2))$, the extractor's output will be close to uniform over its range, which is of size $N^{\Theta(\log(N)/\sigma^2)}$. In the YES case, where $D$'s Shannon entropy was $\log(S)$, the extractor's output will have a support that is significantly smaller than its range. The verifier flips an unbiased coin and, depending on its coin toss, sends to the prover a sample either from the extractor's output or from the uniform distribution over the range. The prover should guess the verifier's coin flip. Soundness follows because in the NO case the two distributions are statistically close. Completeness also follows, because in the YES case the distributions are far. However, the honest prover's running time is huge, as it needs to check whether the sample it received is in the support of the extractor's output, and this requires

10

runtime $N^{\Theta(\log(N)/\sigma^2)}$.

## 2.3 Organization of The Paper

Section 3 contains definition and preliminaries. Sections 4 and 5 comprise the two main pillars of our main result: the first of the two contains a doubly-efficient proof system for obtaining a *tagged sample* of distribution $D$, namely, a collection of samples drawn from $D$, and a set of claims about the probability under $D$ of each sample. The verification of this set of claims is performed in two steps, and includes the *collisions matching tester* constructed in Section 5. Finally, in Section 6 we put all these components together, and construct a doubly-efficient histogram reconstruction protocol, which we then leverage for doubly efficient verification of label-invariant properties.

# 3 Preliminaries

## 3.1 Distributions - General Definitions

Without loss of generality, and for the sake of simplicity of notation ahead, we consider all finite domains to be subsets of $\mathbb{N}$.

**Definition 3.1.** *The statistical distance between distributions $P$ and $Q$ over a finite domain $X$ is defined as:*

$$\Delta_{SD}(P,Q) = \frac{1}{2} \sum_{x \in X} |P(x) - Q(x)|$$

**Claim 3.2.** *Let $P, Q$ be distributions over a domain $\mathcal{X}$ such that $\Delta_{SD}(P,Q) = \delta$. Then:*

$$\max_{A \subseteq \mathcal{X}} (P(A) - Q(A)) = \delta$$

*Proof.* Define $A = \{x \in \mathcal{X} \ : \ P(x) > Q(x)\}$. Observe that by definition:

$$\Delta_{\mathrm{SD}}(P,Q) = \frac{1}{2} \sum_{x \in \mathcal{X}} |P(x) - Q(x)| \tag{5}$$

$$= \frac{1}{2} \sum_{x \in A} (P(x) - Q(x)) + \frac{1}{2} \sum_{x \in \mathcal{X} \setminus A} (Q(x) - P(x)) \tag{6}$$

$$= \frac{1}{2} (P(A) - Q(A)) + \frac{1}{2} (Q(\mathcal{X} \setminus A) - P(\mathcal{X} \setminus A)) \tag{7}$$

$$= \frac{1}{2} (P(A) - Q(A)) + \frac{1}{2} (1 - Q(A) - (1 - P(A))) \tag{8}$$

$$= P(A) - Q(A) \tag{9}$$

Moreover, since by definition for every $x \in A$ it holds that $P(x) > Q(x)$, then:

$$A \in \arg\max_{X \subseteq \mathcal{X}} (P(X) - Q(X))$$

As taking any element out of $A$ will decrease the value of $P(A) - Q(A)$, and any element added to $A$ has to be from $\{x \in \mathcal{X} : Q(x) \geq P(x)\}$, and as such, it won't increase $P(A) - Q(A)$. $\qquad \square$

**Definition 3.3** (($N, \xi$)-bucket of a distribution)**.** *The $\ell$'th $(N, \xi)$ bucket of distribution $P$ is:*

$$B_\ell^P = \left\{ x \in Supp(P) : P(x) \in \left[ \frac{e^{\ell\xi}}{N}, \frac{e^{(\ell+1)\xi}}{N} \right) \right\}$$

**Definition 3.4** (($N, \xi$)-histogram of a distribution)**.** *The $(N, \xi)$-histogram of distribution $P$ is the collection $\{p_\ell\}_{\ell\in\mathbb{Z}}$. We also consider the histogram with all elements of mass lighter than $\frac{\xi^2}{N}$ collected to one bucket, $B_L^P$. That is:*

$$B_L^P = \left\{ x \in Supp(P) : P(x) \leq \frac{\xi^2}{N} \right\}, \quad p_L = P\left( B_L^P \right)$$

This is motivated by two reasons: for a distribution $P$ over domain $[N]$, $p_L \cdot N \leq \xi^2$, and so, it holds that $B_L^P$ accounts for at most $\xi^2$ mass (which we consider small). Therefore, when we consider the $(N, \tau)$-histogram of some distribution, we in fact don't consider the bucket indices to be taken from $\mathbb{Z}$, but from a smaller set:

**Definition 3.5** (Number of bucket buck$(N, \xi)$)**.** *Given parameters $N$ and $\xi$, we consider $b(N, \xi)$ to be the number of buckets with all buckets $\ell$ for which $\frac{e^{\ell\xi}}{N} \geq \frac{\xi^2}{N}$ are collected into one.*

*Observe that $b(N, \xi) = \left\lceil \frac{\log N}{\xi} \right\rceil + \left\lceil \frac{2\log 1/\xi}{\xi} \right\rceil = O\left( \log N/\xi \right)$*

**Convention 3.6.** *In the paper, unless stated otherwise, whenever we talk about bucket indices we consider them to be taken from the index set:*

$$\mathcal{I}_{(N,\xi)} = \left\{ -\left\lceil \frac{2\log 1/\xi}{\xi} \right\rceil, \ldots, -1, 0, 1, \ldots, \left\lceil \frac{\log N}{\xi} \right\rceil \right\}$$

## 3.2   Relabeling Distance

The reader is referred to [HR22] for proofs for all claims in this section.

**Definition 3.7** (Permutation of a distribution)**.** *For a distribution $P$ over a domain $\mathcal{X}$, and a permutation $\pi$ over the same domain, we define $\pi(P)$ as the distribution that satisfies for every $x \in \mathcal{X}$: $\pi(P)(x) = P(\pi^{-1}(x))$.*

**Definition 3.8.** *For any set $A$, perm$(A)$ is the set of all permutations over the set $A$.*

**Definition 3.9** (Relabeling distance)**.** *Let $P$ and $Q$ be distributions over finite domains $\mathcal{X} \subseteq \mathbb{N}$, and $\mathcal{Y} \subseteq \mathbb{N}$ respectively. The relabeling distance between $P$ and $Q$ is defined to be:*

$$\Delta_{RL}(P, Q) = \min \left\{ \Delta_{SD}(P, \pi(Q)) \ : \pi \in perm\left( \mathbb{N} \right) \right\}$$

**Claim 3.10.** *Let $P, Q, R$ be any three distributions over finite domains $\mathcal{X}, \mathcal{Y}$, and $\mathcal{Z}$ respectively. The Relabeling Distance satisfies:*

- $\Delta_{RL}(P, Q) \geq 0$, *and* $\Delta_{RL}(P, Q) = 0$ *iff there exists a permutation* $\sigma \in perm(\mathbb{N})$ *such that* $P = \sigma(Q)$.

- *Symmetry:* $\Delta_{RL}(P, Q) = \Delta_{RL}(Q, P)$.

12

- *Triangle inequality:* $\Delta_{RL}(P, R) \leq \Delta_{RL}(P, Q) + \Delta_{RL}(Q, R)$

**Definition 3.11.** *For any finite distribution $Q$, and any $(N, \xi)$-histogram $\{p_j\}_j$, define:*

$$\Delta_{RL}(Q, \{p_j\}_j = \min_{P \text{ has histogram } \{p_j\}_j} \Delta_{SD}(Q, P)$$

*This definition also extends to the distance between two histograms.*

**Claim 3.12.** *Let $\{q_j\}_j$ and $\{p_j\}_j$ be two $(N, \xi)$-histograms. For every $\varepsilon \geq 0$, if $\frac{1}{2} \sum_j |p_j - q_j| \leq \varepsilon$, then,*

$$\Delta_{RL}(\{q_j\}_j, \{p_j\}_j) \leq e^{\xi}\varepsilon + e^{\xi}(e^{\xi} - 1)$$

**Claim 3.13.** *For any two distributions $P, Q$ over the domain $[N]$. Let $\pi^{ord} : [N] \to [N]$ be the permutation that satisfies the property for every $i, j \in [N]$, if $P(i) < P(j)$ then $\left(\pi^{ord}(Q)\right)(i) \leq \left(\pi^{ord}(Q)\right)(j)$. Then, it holds that:*

$$\Delta_{RL}(P, Q) = \Delta_{SD}(P, \pi^{ord}(Q))$$

**Proposition 3.14** (Histogram distance estimator). *For every $\xi \leq 0.1$ there exists an algorithm that runs in $O\left(\log(N)/\xi\right)$ time and given parameters $N$, $\xi$, as well as two $(N, \xi)$-histograms $\{p_j\}_j$ and $\{q_j\}_j$, outputs $d$ such that $|d - \Delta_{RL}\left(\{p_j\}_j, \{q_j\}_j\right)| \leq 7\xi$.*

## 3.3 Testing and Verifying Distribution Properties

**Theorem 3.15** (Folklore distribution learner [Gol17]). *There exists an algorithm that given sample access to a distribution $P$ over the domain $[N]$, and an accuracy parameter $\alpha \in (0, 1)$, it runs in time $\tilde{O}(N/\alpha^2)$, takes $O(N/\alpha^2)$ samples, and with probability at least 0.99 outputs a full description of a distribution $P_{approx}$ such that $\Delta_{SD}(P, P_{approx}) \leq \alpha$.*

**Definition 3.16** (Distribution tester for property $\Pi$). *Let $\delta$ be some distance measure between distributions, and $\Pi$ be some collection of finite distributions. Denote $\Pi_N = \Pi \cap \Delta_N$ (where $\Delta_N$ is the set of all distributions over domain of size at most $N$). A tester $T$ of property $\Pi$ is a probabilistic oracle machine, that on input parameters $N$ and $\varepsilon$, and oracle access to a sampling device for a distribution $D$ over a domain of size $N$, outputs a binary verdict that satisfies the following two conditions:*

1. *If $D \in \Pi_N$, then $\Pr(T^D(N, \varepsilon) = 1) \geq 2/3$.*

2. *If $\delta(D, \Pi_N) > \varepsilon$, then $\Pr(T^D(N, \varepsilon) = 0) \geq 2/3$.*

In the context of this work, the relevant distance measure is *statistical distance* as defined above. An extension of this definition, introduced by Parnas, Ron, and Rubinfeld [PRR06] is the following:

**Definition 3.17** (($\varepsilon_c, \varepsilon_f$)-tolerant distribution property tester). *For parameters $\varepsilon_c, \varepsilon_f \in [0, 1]$ such that $\varepsilon_c < \varepsilon_f$, a $(\varepsilon_c, \varepsilon_f)$-tolerant tester $T$ of property $\Pi$ is a probabilistic oracle machine, that on inputs $N, \varepsilon_c, \varepsilon_f$ and given oracle access to a sampling device for distribution $D$ over a domain of size $N$, outputs a binary verdict that satisfies the following two conditions:*

1. *If $\delta(D, \Pi_N) \leq \varepsilon_c$, then $\Pr(T^D(N, \varepsilon_c, \varepsilon_f) = 1) \geq 2/3$.*

2. If $\delta(D, \Pi_N) \geq \varepsilon_f$, then $\Pr(T^D(N, \varepsilon_c, \varepsilon_f) = 0) \geq 2/3$.

Note that a tolerant distribution test is for some property $\Pi$ is at least as hard as a standard non-tolerant tester for the same property.

Our main result is a double efficient interactive proof system for many tolerant testing problems. The following definition sets the framework for this work. It is based on the setting presented in the seminal work of Goldwasser, Micali, and Rackoff [GMR85], and it is an extension of the definition presented by Chiesa and Gur [CG18] that includes tolerant testing, that seeks to have small honest prover complexity, in the vein of [GKR15].

**Definition 3.18** (Proof system for tolerant distribution testing problems). *A proof system for a tolerant distribution testing problem $\Pi$ with parameters $\varepsilon_c$ and $\varepsilon_f$ is a two-party game, between a verifier executing a probabilistic polynomial time strategy $V$, and a prover that executes a strategy $P$. Given that both $V$ and $P$ have black-box sample access to distribution $D$ over the domain $[N]$, and are given $N$, the interaction should satisfy the following conditions:*

- **Completeness:** *For every $D$ over domain of size at most $N$, such that $\Delta_{SD}(D, \Pi_N) \leq \varepsilon_c$, the verifier $V$, after interacting with the prover $P$, accepts with probability at least $2/3$.*

- **Soundness:** *For every $D$ over domain of size at most $N$ such that $\Delta_{SD}(D, \Pi_N) \geq \varepsilon_f$, and every cheating strategy $P^*$, the verifier $V$, after interacting with the prover $P^*$, rejects with probability at least $2/3$.*

*The complexity measures associated with the protocol are: the sample complexity of the verifier as as the honest prover (strategy $P$), the communication complexity, the runtime of both agents, and the round complexity (how many messages were exchanged).*

**Definition 3.19** (Label invariant distribution property). *A distribution property $\Pi$ is called label invariant if for all $N \in \mathbb{N}$, it holds that any permutation $\sigma$ over $N$ elements satisfies that $D \in \Pi_N$ if and only if $\sigma(D) \in \Pi_N$.*

# 4 Bucket Size Verification Protocol

## 4.1 Protocol Overview

In this section, given sample access to a distribution $D$ over domain $[N]$, we construct a protocol for obtaining a collection of claims about the probability under $D$ of each element in a set $S \subseteq [N]$ of size roughly $\sqrt{N}$.

Concretely, given some accuracy parameter $\tau \in (0, 0.1)$, at the end of the protocol, the verifier is left with $(z_i, \pi(z_i))_{i \in [s]}$ where $s = \widetilde{\theta}\left(\sqrt{N} \cdot \mathsf{poly}(\tau^{-1})\right)$, $(z_i)_{i \in [s]}$ was drawn i.i.d. by $D$, and for every $i \in [s]$, it is alleged that $\pi(z_i) \approx D(z_i)$.

If the prover is honest, then indeed $\pi(z_i) = (1 \pm O(\tau)) \cdot D(z_i)$, and a honest prover strategy with input the same as that of the verifier (black-box sample access to $D$, parameters $N$ and $\tau$) can be implemented in time roughly linear in $N$. If the prover is dishonest (and potentially computationally unbounded), however, we characterize the distance between the alleged probability under $D$ of each element $(\pi(z_i))_{i \in [s]}$ to their true probability $(D(z_i))_{i \in [s]}$ through the following variables:

$$x_{\ell, j} = \frac{\left| \left\{ i \in [s] : z_i \in B_\ell^D, \pi(z_i) \in \left[ \frac{e^{j\tau}}{N}, \frac{e^{(j+1)\tau}}{N} \right] \right\} \right|}{\left| \left\{ i \in [s] : z_i \in B_\ell^D \right\} \right|}$$

Recall that $B_\ell^D = \left\{ x \in [N] : D(x) \in \left[ \frac{e^{\ell\tau}}{N}, \frac{e^{(\ell+1)\tau}}{N} \right) \right\}$. That is, $x_{\ell, j}$ should be thought of as the fraction of elements with mass approximately $e^{\ell\tau}/N$ that were claimed to have mass approximately $e^{j\tau}/N$ at the end of the protocol (indeed, if the prover is honest, then $x_{\ell, j} = \mathbb{1}_{\ell=j}$). We are guaranteed that if the verifier accepted, no matter how the prover chose to respond or cheat, the following holds:

**The Soundness Guarantee.** *Fix a sample and set of claims $(z_i, \pi(z_i))_{i \in [s]}$ obtained through the protocol. Consider the following variables induced by these claims, for every $\ell \in \mathbb{Z}$:*

$$v_\ell = \frac{\left| \left\{ i \in [s] : \pi(z_i) \in \left[ \frac{e^{\ell\tau}}{N}, \frac{e^{(\ell+1)\tau}}{N} \right] \right\} \right|}{s}, \quad \widehat{q}_\ell = \frac{\left| \left\{ i \in [s] : D(z_i) \in \left[ \frac{e^{\ell\tau}}{N}, \frac{e^{(\ell+1)\tau}}{N} \right] \right\} \right|}{s}$$

*That is, $\widehat{q}_\ell$ is the empirical mass of $B_\ell^D$ according to the sample $(z_i)_{i \in [s]}$, and thus, it is strongly concentrated around $D\left(B_\ell^D\right) = q_\ell$, while $v_\ell$ is the alleged empirical mass of the same set, and can potentially be far from $\widehat{q}_\ell$. In particular by definition, $v_j = \sum_\ell \widehat{q}_\ell x_{\ell, j}$. At the end of the protocol, with high probability over the verifier's coin tosses and samples, if they accepted, then it must be that:*

$$N v_j e^{-j\tau} \approx \sum_\ell N \widehat{q}_\ell e^{-\ell\tau} x_{\ell, j} \tag{10}$$

**Digest of soundness guarantee.** The soundness condition claims that any interaction at the end of which Inequality (10) doesn't hold, will be rejected with high probability. I.e. by definition it must hold that $v_j = \sum_\ell \widehat{q}_\ell x_{\ell, j}$, and the protocol enforces another condition over $\{x_{\ell, j}\}_{\ell, j}$, namely $v_j e^{-j\tau} \approx \sum_\ell \widehat{q}_\ell x_{\ell, j} e^{-\ell\tau}$. This can be interpreted as a condition over the alleged size of each bucket: recall that allegedly, for every $\ell$, $v_\ell = \widehat{q}_\ell$, and since by definition, $\left| B_\ell^D \right| \in \left[ \frac{q_\ell}{e^{(\ell+1)}/N}, \frac{q_\ell}{e^\ell/N} \right] \approx$

$\left[\frac{\widehat{q_\ell}}{e^{(\ell+1)}/N}, \frac{\widehat{q_\ell}}{e^{\ell}/N}\right]$, the verifier considers the left-hand side of Approximate Equality (10), $\frac{v_j}{e^{j\tau}/N}$, as the *alleged approximate size* of $B_j^D$ according to $(z_i, \pi(z_i))_{i\in[s]}$. Thus, we think of the soundness condition of the protocol as requiring the prover to lie in such a way that *the alleged size of each bucket* conforms to the above approximate equality.

Note that if the prover behaves dishonestly, the verifier in this protocol is not required to reject with high probability. Indeed, catching a cheating prover is done in several stages, and the protocol presented in this section is only one step in this process. In order to complete the verification of the above-mentioned claims regarding the alleged probability of the elements sampled, we require the tester described in Section 5. The way both these procedures are combined can be found in Section 6. We continue to outline how the *Bucket Size Verification Protocol* works:

**Step I: obtaining tagged samples.** The verifier flips a balanced roughly coin $2s$ times, and obtains $b \in \{0,1\}^{2s}$. Every time the coin shows 0, the verifier draws a sample from $D$, and otherwise it draws a sample from $U_{[N]}$. We think of the samples drawn $(z_1, z_2, \ldots, z_{2s})$ as actually being composed of two different interweaved samples, and denote by $S_0 = \{i : b_i = 0\}$ the samples drawn by $D$, and $S_1 = [2s] \setminus S_0$, the samples drawn by $U_{[N]}$. For simplicity assume both sets are of size $s$. The sample $(z_1, \ldots, z_{2s})$ is sent to the prover (note that the choice of coins $b$ is not revealed to the prover).

For every $i \in [2s]$, the prover replies with $\pi(z_i)$ (often referred to as *tags* in the paper, as they are thought of placement in buckets) as explained above, or $\perp$ if the sample was drawn from $U_{[N]}$ outside the support of $D$. If the prover is honest, the tags are all correct.

**Step II: bucket size consistency test.** The verifier computes the empirical $(N, \tau)$-histogram of $D$ implied by the prover's answer. For every bucket index $j$, it considers the set of elements drawn from $D$ tagged as belonging to bucket $j$:

$$S_0^j = \left\{ i \in S_0 : \pi(z_i) \in \left[\frac{e^{j\tau}}{N}, \frac{e^{(j+1)\tau}}{N}\right) \right\}$$

Then, it sets $v_j = \frac{|S_0^j|}{|S_0|}$. If the prover is honest, we expect this quantity to be very strongly concentrated around the true mass of the bucket $q_j = D(B_j^D)$. Since, for every bucket $j$ it holds that $\left|B_j^D\right| \in \left(\frac{q_j}{e^{(j+1)\tau}/N}, \frac{q_j}{e^{j\tau}/N}\right]$, the verifier derives from $v_j$ an alleged size for the $j$'th bucket: $\frac{v_j}{e^{j\tau}/N}$.

A main idea behind the bucket size consistency check is that the verifier actually gets *two* approximations for the size of each bucket, and checks that they are close. One is obtained through considering the samples in $S_0$, as explained above, and the other, through $S_1$. Focusing on the second approximation, for every bucket index $j$, the verifier considers the following quantity:

$$w_j = \frac{\left|\left\{ i \in S_1 : \pi(z_i) \in \left[\frac{e^{j\tau}}{N}, \frac{e^{(j+1)\tau}}{N}\right) \right\}\right|}{s_1}$$

As before, the meaning of this quantity in the case that the prover is honest is the fraction of samples drawn from $U_{[N]}$ that landed in the $j$'th bucket of $D$. Therefore if the prover is honest, we expect $w_j$ to be (strongly) concentrated around $\frac{|B_j^D|}{N}$. And so, we can think of $(N \cdot w_j)$ as another

approximation of $\left|B_j^D\right|$. The verifier then accepts when these two approximations are indeed close, that is, if the following holds:

$$w_j \approx v_j e^{-j\tau}$$

After passing this test, the verifier outputs $(z_i, \pi(z_i))_{i \in S_0}$.

**Proof intuition.** The completeness of the protocol follows immediately from the explanations above. In short, we expect both $\frac{N v_j}{e^{j\tau}}$ as well as $(N \cdot w_j)$ to be concentrated around $\left|B_j^D\right|$, and thus close. The soundness condition is a bit more involved. For every two bucket indices $\ell$ and $j$, we consider the variables $x_{\ell,j}, y_{\ell,j} \in [0,1]$ as a characterization of *how* the prover cheated. We explained above the definition of $\{x_{\ell,j}\}_{\ell,j}$ as a characterization of how the prover cheated over $S_0$. The collection $\{y_{\ell,j}\}_{\ell,j}$ is defined analogously with respect to $S_1$ (see Definition 4.2 and following discussion for more detail). We capitalize strongly on the fact that for (almost) every $\ell$ and $j$, $x_{\ell,j}$ and $y_{\ell,j}$ are closely related. Intuitively, this can be explained through the following important observation:

**Observation 4.1.** *A cheating prover that wants to mistag the samples that were drawn according to $D$ (e.g. tag half the samples drawn from $D$ that landed in bucket $B_\ell^D$ as if they belong to bucket $j$) must also, in the process, mistag samples that were drawn according to $U_{[N]}$ in a similar pattern (following the above example, half of the samples drawn according to $U_{[N]}$ that landed in $B_\ell^D$, will be tagged as belonging to bucket $j$ according to $D$).*

*This is justified by the fact that the prover doesn't know $b \in \{0,1\}^s$, and so, when considering the set of samples in $(z_1, \ldots, z_s)$ that truly landed in the set $B_\ell^D$, it is unable to determine which was drawn from $D$ and which from $U_{[N]}$.*

And so, if we write $m_\ell = \frac{|i \in S_1 : z_i \in B_\ell^D|}{|S_1|}$ (i.e. the empirical mass of the set $B_\ell^D$ according to the uniform distribution over the domain, and the samples in $S_1$), we can rewrite $w_j = \sum_\ell m_\ell y_{\ell,j}$. Now, from the same reasoning as above, $m_\ell$ is strongly concentrated around $\left|B_\ell^D\right|/N$, which is roughly $N q_\ell e^{-\ell\tau}/N \approx \widehat{q}_\ell e^{-\ell\tau}$. And so, plugging $m_\ell \approx \widehat{q}_\ell e^{-\ell\tau}$ and $x_{\ell,j} \approx y_{\ell,j}$, we get with high probability:

$$w_j \approx \sum_\ell \widehat{q}_\ell e^{-\ell\tau} x_{\ell,j}$$

We thus get that by verifying that $w_j \approx v_j e^{-j\tau}$, we get the desired result, namely:

$$N v_j e^{-j\tau} \approx N \sum_\ell \widehat{q}_\ell e^{-\ell\tau} x_{\ell,j} \tag{11}$$

**Honest prover runtime.** A prover that approximates the probability of each element with mass at least $\frac{\tau}{N}$ to a multiplicative factor of $\tau/10$ can be implemented in time $\widetilde{O}(N)\mathsf{poly}(\tau^{-1})$. Such approximation might not put every element in its correct bucket, as elements with probability close to the margins of the buckets might be placed in adjacent buckets. This does not affect the completeness of the protocol, which can withstand such errors. See Remark 4.14 for more detail.

**Technical remarks.** The above intuition disregards issues of measure concentration, and the treatment of *small buckets* for which we cannot ensure that $x_{\ell,j} \approx y_{\ell,j}$. Consequently, taking these points into account, the true soundness guarantee is a lower bound on $N v_j e^{-j\tau}$, rather than an approximate equality.

## 4.2 Bucket Size Verification Protocol

---

**Protocol 4.1.1: Bucket Size Verification Protocol**

**Input:** parameters $N \in \mathbb{N}$, $\tau \in (0,1)$, as well as sample access to distribution $D$ over domain $[N]$. $D$ is assumed to satisfy $\forall x \; D(x) \leq \frac{\tau}{\sqrt{N}}$

1. V: toss a balanced coin $2s$ times for $s = \sqrt{N} \cdot \mathsf{poly}(\log N, \tau^{-1})$ to obtain $b \in \{0,1\}^{2s}$. For every $i \in [2s]$, if $b_i = 0$, draw $z_i \leftarrow D$; otherwise, draw $z_i \leftarrow U_{[N]}$. Send $(z_1, z_2, \ldots, z_s)$ to P.

2. P: for every $i \in [2s]$: if $z_i \in \mathrm{Supp}(D)$, set $\pi(z_i)$ so that $\pi(z_i) \in D(z_i)[1 - \tau/10, 1 + \tau/10]$ (see Algorithm 4.12.1 for details), and otherwise set $\pi(z_i) = \perp$.

3. V: set $S_0 = \{i : b_i = 0\}$, and $|S_0| = s_0$, as well as $S_1 = \{i : b_i = 1\}$, with $|S_1| = s_1$. If $|s_0 - s/2| > s/6$, reject. For every $i \in [2s]$ such that $\pi(z_i) \neq \perp$, also set $\mathrm{tag}(z_i) = \lfloor \log\left(N \cdot \pi(z_i)\right)/\tau \rfloor$. Do the following:

   (a) **Basic consistency test.** Check that for all $i_1, i_2 \in [2s]$, if $z_{i_1} = z_{i_2}$, $\mathrm{tag}\left(z_{i_1}\right) = \mathrm{tag}\left(z_{i_2}\right)$, and that for all $i \in S_0$, $\pi(z_i) \neq \perp$.

   (b) **Verify bucket sizes.** For every $j$, define $v_j = \frac{|\{i \in S_0 : \mathrm{tag}(z_i) = j\}|}{s_0}$, and $w_j = \frac{|\{i \in S_1 : \mathrm{tag}(z_i) = j\}|}{s_1}$.
   Reject if there exists such $j$ such that $\frac{e^{j\tau}}{N} \geq \frac{\tau}{N}$ for which $v_j \geq \frac{\tau^2}{\log N}$ and:
   $$\left| w_j - v_j e^{-j\tau} \right| \geq 4\tau v_j e^{-j\tau} \tag{12}$$

4. V: output $(z_i, \pi(z_i))_{i \in S_0}$.

---

In order to formally claim the completeness and soundness conditions that apply to Protocol 4.1.1, we first define a set of variables that characterize all possible prover responses in the protocol:

**Definition 4.2.** *Any prover response in Protocol 4.1.1 induces the following collections of variables. For every two buckets $\ell$ and $j$ define:*

$$x_{\ell,j} = \frac{\left|\left\{i \in S_0 : z_i \in B_\ell^D, \; tag(z_i) = j\right\}\right|}{\left|\left\{i \in S_0 : z_i \in B_\ell^D\right\}\right|}$$

$$y_{\ell,j} = \frac{\left|\left\{i \in S_1 : z_i \in B_\ell^D, \; tag(z_i) = j\right\}\right|}{\left|\left\{i \in S_1 : z_i \in B_\ell^D\right\}\right|}$$

To give some intuition as to meaning of these variables, as expained above in the protocol overview, and as will be further explained in the analysis, consider some prover strategy: the prover receives a collection of samples from the verifier, some drawn by $D$ and some by $U_{[N]}$. The first set of variables $\{x_{\ell,j}\}_{\ell,j}$ characterizes (cheating) prover behavior over samples drawn by the verifier according to $D$, and $\{y_{\ell,j}\}_{\ell,j}$ captures (cheating) behavior over samples drawn according to $U_{[N]}$. It is important to note that the prover doesn't know for each sample whether it was drawn according to $D$ or according to $U_{[N]}$. More concretely, the variable $x_{\ell,j}$ (respectively $y_{\ell,j}$) describes the fraction of samples drawn according from $D$ ($U_{[N]}$) that landed in bucket $B_\ell^D$, but were reported as belonging to $D$-bucket $j$. In particular, if the prover is honest we get $y_{\ell,j} = x_{\ell,j} = 1$ if $\ell = j$ and $y_{\ell,j} = x_{\ell,j} = 0$ otherwise.

**Lemma 4.3.** *Protocol 4.1.1 satisfies the following conditions: the verifier runs in time $\widetilde{O}(s)$, and draws $O(s)$ samples, for $s = \widetilde{O}\left(\sqrt{N}\right)\mathsf{poly}(\tau^{-1})$; the honest prover, given the same input as the verifier, can be implemented with both its runtime and sample complexity be of magnitude $\widetilde{O}(N)\,\mathsf{poly}\left(\tau^{-1}\right)$. At the end of the interaction the verifier either rejects or outputs $(z_i, \pi(z_i))_{i\in[s]}$, such that the following conditions hold:*

- **Completeness.** *If the prover is honest, then, with probability at least $0.95$ over the samples and coin-tosses of both the verifier and the prover, the verifier accepts and for every $i \in [s]$, $\pi(z_i)$ is a $\tau$-approximation of $D(z_i)$.*

- **Soundness.** *No matter what strategy a cheating prover might employ, with probability at least $0.95$ over the samples and the coin tosses of the verifier, either the verifier rejects, or for every $j$ such that $\frac{e^{j\tau}}{N} \geq \frac{\tau}{N}$ and $v_j \geq \frac{\tau^2}{\log N}$:*

$$v_j e^{-j\tau} \geq (1 - 18\tau) \sum_{\ell:\frac{e^{\ell\tau}}{N} \geq \frac{\tau^2}{N}} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau} - \mathsf{poly}\left(\log N, \tau^{-1}\right) \cdot \frac{1}{s} \tag{13}$$

## 4.3  Proof of Lemma 4.3

### 4.3.1  Completeness of Protocol 4.1.1

For simplicity, we analyze the protocol's completeness under the simplifying assumption that the honest prover tags every sample with it's true bucket. Looking ahead, since we want a doubly-efficient honest prover, the tags will be according to multiplicative approximations to the true probabilities, and some samples might be slightly mis-tagged. This does not have much of an effect on the completeness analysis, see Remark 4.14 for further details.

With high probability the number of samples drawn according to $D$ and the number of samples drawn according to $U_{[N]}$ will be roughly $s$. Formally:

**Claim 4.4.** *For every $s > 500\tau^{-2}$, with probability at least $0.999$, $\left|\sum_{i\in[s]} b_i - s\right| < s \cdot \tau$*

*Proof.* By Hoeffding's Inequality, for $s > 500\tau^{-2}$:

$$\Pr_b\left(\left|\sum_{i\in[2s]} b_i - s\right| > s\tau\right) < e^{-2s\cdot\tau^2/80} < 0.001$$

$\square$

We choose $s > 500\tau^{-2}$, so the verifier is unlikely to reject due to a bad choice of bits $(b_i)_i$.

Next, observe that by definition, a honest prover passes the *basic consistency test* (Step (3a)), and provides $(\pi(z_i))_{i\in S_0}$ such that $\pi(z_i) \in D(z_i)[1 - \tau/10, 1 + \tau/10]$. And so, in order to show completeness, we need to show that for every bucket index $j$ such that $v_j \geq \frac{\tau^2}{\log N}$, Inequality (12) holds.

Following the line of argument in the protocol overview in Section 4.1, we show that if the prover is honest, then for every $j$ such that $v_j \geq \frac{\tau^2}{\log N}$, it holds that $v_j$ and $w_j$ are tightly concentrated

19

around $q_j = D(B_j^D)$ and $\frac{|B_j^D|}{N}$ respectively. The claim again follows from the multiplicative Chernoff bound, but since we will come back to the connection between the empirical mass of buckets and their true mass according to the distribution from which the sample was drawn, we introduce the following definition and claim:

**Definition 4.5** ( $(\mathcal{Z}, R, \tau)$-characteristic sample)**.** *Let $R \geq 1$ and $\tau > 0$ be positive real numbers, $\mathcal{Z}$ a partition of the domain $[N]$, of size $|\mathcal{Z}|$, then a sample $S = (z_1, \ldots, z_t)$ drawn i.i.d. according to distribution $P$ over $[N]$ is $(\mathcal{Z}, R, \tau)$-characteristic if for every set $A \in \mathcal{Z}$, such that $P(A) > \frac{1}{R}$ it holds that:*

$$\left| \frac{|\{i \in [t] : z_i \in A\}|}{t} - P(A) \right| < \tau \cdot P(A)$$

**Claim 4.6.** *For domain $[N]$, let $R, \mathcal{Z}, P$ be as specified in Definition 4.5. Let $S$ be a sample of size $t \geq \frac{R \log(1000|\mathcal{Z}|)}{\tau^2}$ drawn i.i.d. by the distribution $P$. Then, with probability at least 0.999, $S$ is $(\mathcal{Z}, R, \tau)$-characteristic with respect to $P$.*

*Proof.* Fix $A \in \mathcal{Z}$. For every $i \in [t]$, define the indicator that $S_i \in A$ to be $\mathbb{1}_{z_i \in A}$. Note that $\frac{|\{i \in [t] : z_i \in A\}|}{s} = \frac{1}{s} \sum_{i \in [t]} \mathbb{1}_{S_i \in A}$, and that this sum is composed of independent variables. Since for every $i$, $\mathbb{E}_{z_i \sim P} [\mathbb{1}_{z_i \in A}] = P(A)$, and all the samples were drawn i.i.d., $\mathbb{E}_{z \sim P^s} \left[ \frac{1}{t} \sum_{i \in [t]} \mathbb{1}_{z_i \in A} \right] = P(A)$. By applying the multiplicative Chernoff bound with $t \geq \frac{R \log(1000|\mathcal{Z}|)}{\tau^2}$ we conclude that for every $A \in \mathcal{Z}$ such that $P(A) > \frac{1}{R}$, with probability at most $\frac{1}{1000|\mathcal{Z}|}$:

$$\left| \frac{|\{i \in [t] : z_i \in A\}|}{t} - P(A) \right| \geq \tau \cdot P(A)$$

Taking the union bound over $A \in \mathcal{Z}$ provides the desired result. $\square$

We can now think of $\{B_\ell^D\}_\ell$ as a partition of $[N]$ of size $b(N, \tau)$. Since the sample is considerably larger than $b(N, \tau)$, it follows that with very high probability the sample $V$ draws will be characteristic with respect to this partition. More concretely, we define the following variables that characterize any sample drawn, regardless of the prover's response:

**Definition 4.7.** *The (true) fraction of samples drawn according to $U_{[N]}$ that landed in the set $B_\ell^D$ is:*

$$m_\ell = \frac{\left| \{i \in S_1 : z_i \in B_\ell^D\} \right|}{s_1}$$

*The (true) fraction of samples drawn according to $Q$ that landed in the set $B_\ell^D$ is:*

$$\widehat{q}_\ell = \frac{\left| \{i \in S_0 : z_i \in B_\ell^D\} \right|}{s_0}$$

Following Definition 4.2, for every prover response:

$$w_j = \sum_\ell m_\ell \cdot y_{\ell, j}$$

$$v_j = \sum_\ell \widehat{q}_\ell \cdot x_{\ell, j}$$

In this section we focus on the case that the prover is honest, and for every bucket index $j$, we assume $v_j = \widehat{q}_j$, as well as $w_j = m_j$.

**Corollary 4.8.** *With probability at least* $0.998$, *for every* $j$ *such that* $q_j > \frac{\log(1000b(N,\tau))}{s\tau^2}$: $|\widehat{q}_j - q_j| < \tau \cdot q_j$

*Proof.* By Claim 4.4 and choice of $s$, with probability at least $0.999$, $s_0 > s/3$. By Claim 4.3.1 and the choice of $s$, give that $s_0 > s/3$, it holds that with probability at least $0.999$, the sample $(z_i)_{i \in S_0}$ is $\left(\{B_\ell^D\}_\ell, \frac{s\tau^2}{\log(1000b(N,\tau))}, \tau\right)$-characteristic with respect to $Q$, and in particular, if $q_j > \frac{\log(1000b(N,\tau))}{s\tau^2}$, $|v_j - q_j| = |\widehat{q}_j - q_j| < \tau \cdot q_j$. $\qquad\square$

Since the verifier has only access to $v_j$ and doesn't know $q_j$, we wish to have a guarantee that $v_j$ and $q_j$ are close, given that $v_j$ is large enough (and not given that $q_j$ is large enough, as the previous claim shows).

**Claim 4.9.** *If the prover is honest, then with probability at least* $0.96$ *over the choice of* $(b_i)_{i \in [s]}$ *and* $(z_1, \ldots, z_s)$, *for every* $j$ *such that* $v_j > \frac{\tau^2}{\log N}$: $|v_j - q_j| < \tau \cdot q_j$

*Proof.* Assume $s_0 > s/3$. Since the prover is honest, for every bucket index $j$ it holds that $v_j = \widehat{q}_j$, as well as $\mathbb{E}_{(z_i)_{i \in S_0}}[v_j] = q_j$. By Markov's Inequality, for all $j$, with probability at most $\frac{1}{1000b(N,\tau)}$, $v_j \geq 1000b(N,\tau)q_j$. In particular, with probability at least $1 - \frac{1}{1000b(N,\tau)}$, every $j$ such that $q_j < \frac{\tau^2}{1000b(N,\tau)\log N}$, it holds that $v_j < \frac{\tau^2}{\log N}$. Taking the union bound over all buckets, with probability at least $0.999$, for every $j$, if $q_j < \frac{\tau^2}{1000b(N,\tau)\log N}$ then $v_j < \frac{\tau^2}{\log N}$. Therefore, if $v_j \geq \frac{\tau^2}{\log N}$, then $q_j > \frac{\tau^2}{1000b(N,\tau)\log N}$, and in particular, through the choice of $s$, $q_j > \frac{\log(1000b(N,\tau))}{s\tau^2}$, and following Corollary 4.8, we conclude with probability at least $0.96$ over the verifier's randomness, for every bucket index $j$ such that $v_j > \frac{\tau^2}{\log N}$, it holds that $|v_j - q_j| < \tau \cdot q_j$. $\qquad\square$

Moving on to show a similar result for the samples drawn from $U_{[N]}$:

**Corollary 4.10.** *With probability at least* $0.99$ *over the choice of* $(z_i)_{i \in [s]}$, *for every* $\ell$ *such that* $\frac{|B_\ell^D|}{N} > \frac{\log(1000b(N,\tau))}{s\tau^2}$:

$$\left|m_\ell - q_\ell e^{-\ell\tau}\right| < 2\tau \cdot q_\ell e^{-\ell\tau}$$

*Proof.* By Claim 4.4, with probability at least $0.999$, $s_1 > s/3$. Assume this is the case. Note that $\frac{|B_\ell^D|}{N} = U_{[N]}\left(B_\ell^D\right)$. By Claim 4.3.1, it holds that with probability at least $0.999$, $(z_i)_{i \in S_1}$ is $\left(\{B_\ell^D\}_\ell, \frac{s\tau^2}{\log(1000b(N,\tau))}, \tau\right)$-characteristic with respect to $U_{[K]}$, i.e. for every $\ell$ such that $\frac{|B_\ell^D|}{N} > \frac{\log(1000b(N,\tau))}{s\tau^2}$:

$$\left|m_\ell - \frac{|B_\ell^D|}{N}\right| < \tau \cdot \frac{|B_\ell^D|}{N}$$

To conclude the proof observe that by definition, $\left|B_\ell^D\right| \in \left[e^{-\tau} \cdot N q_\ell e^{-\ell\tau}, N q_\ell e^{-\ell\tau}\right)$ $\qquad\square$

**Claim 4.11.** *If the prover is honest, then with probability at least* $0.98$ *over the choice of* $(z_i)_{i \in [s]}$ *and* $(b_i)_{i \in [s]}$ *it holds that for all* $\ell$, *if* $v_\ell > \frac{\tau^2}{\log N}$, *then* $\left|w_\ell - q_\ell e^{-\ell\tau}\right| < 2\tau \cdot q_\ell e^{-\ell\tau}$

21

*Proof.* By Claim 4.9, if the prover is honest, then with probability at least 0.96 over the choice of $(z_i)_{i\in[s]}$ and $(b_i)_{i\in[s]}$, it holds that for all $\ell$ such that $v_\ell > \frac{\tau^2}{\log N}$, $|v_\ell - q_\ell| < \tau \cdot q_\ell$, this implies the following:

$$\left|B_\ell^D\right| \geq e^{-\tau} N q_\ell e^{-\ell\tau} \geq e^{-2\tau} N v_j e^{-\ell\tau} \geq e^{-2\tau} \cdot N e^{-\ell\tau} \cdot \frac{\tau^2}{\log N} \geq e^{-2\tau} \frac{\sqrt{N}}{\tau} \cdot \frac{\tau^2}{\log N} \geq \frac{\sqrt{N}\tau}{2\log N}$$

Where the second to last inequality is justified through the assumption that for every $\ell$ such that $B_\ell^D \neq \phi$ it holds that $\frac{e^{\ell\tau}}{N} \leq \frac{\tau}{\sqrt{N}}$. Following the above inequality:

$$\frac{\left|B_\ell^D\right|}{N} \geq \frac{\sqrt{N}\tau}{2N\log N} = \frac{\tau}{\sqrt{N}\log N} \geq \frac{\log\left(1000b(N,\tau)\right)}{s\tau^2}$$

Where the last inequality is justified through the choice of $s$. By so, we conclude that if the prover is honest, for every $\ell$ such that $v_\ell > \frac{\tau^2}{\log N}$, it holds that:

$$\frac{\left|B_\ell^D\right|}{N} \geq \frac{\log\left(1000b(N,\tau)\right)}{s\tau^2}$$

And by Corollary 4.10, with probability at least 0.99 for all such $\ell$:

$$\left|m_\ell - q_\ell e^{-\ell\tau}\right| < 2\tau \cdot q_\ell e^{-\ell\tau}$$

□

We have shown that with high probability if some bucket index $j$ satisfies the condition that $v_j \geq \frac{\tau^2}{\log}$, then $v_j \approx q_j$, or alternatively, $N v_j e^{-j\tau} \approx N q_j e^{-j\tau}$ (Claim 4.9). Similarly, we have shown that with high probability, for every bucket index $j$ for which $\frac{\left|B_\ell^D\right|}{N} > \frac{\log(1000b(N,\tau))}{s\tau^2}$, it also holds that $m_\ell \approx q_\ell e^{-\ell\tau}$ (Claim 4.10). We put both these claims together to show that with high probability the second consistency test passes:

**Claim 4.12.** *If the prover is honest, with probability at least* 0.9, *for all $j$ such that $v_j \geq \frac{\tau^2}{\log N}$ it holds that:*

$$\left|w_j - v_j e^{-j\tau}\right| \leq 4\tau \cdot v_j e^{-j\tau}$$

*Proof.* By Claim 4.9, with probability at least 0.98, for every $j$ such that $v_j > \frac{\tau^2}{\log N}$, $|q_j - v_j| < \tau q_j$. By Claim 4.11, with probability at least 0.98 for every $j$ such that $v_j > \frac{\tau^2}{\log N}$, it holds that $\left|w_j - q_j e^{-j\tau}\right| \leq 2\tau \cdot q_j e^{-j\tau}$. And so, we conclude that with probability at least 0.95 over choice of $(b_i)_{i\in[s]}$ and $(z_i)_{i\in[s]}$, for all $j$ such that $v_j > \frac{\tau^2}{\log N}$:

$$\begin{aligned}
\left|w_j - v_j e^{-j\tau}\right| &\leq \left|w_j - q_j e^{-j\tau}\right| + \left|q_j e^{-j\tau} - v_j e^{-j\tau}\right| \\
&= \left|w_j - q_j e^{-j\tau}\right| + e^{-j\tau}\left|q_j - v_j\right| \\
&\leq 2\tau \cdot q_j e^{-j\tau} + \tau \cdot q_j e^{-j\tau} \\
&\leq 3\tau q_j e^{-j\tau} \\
&< 4\tau v_j e^{-j\tau}
\end{aligned}$$

□

We proved that with high probability, if the prover is honest, the verifier does not reject.

### 4.3.2 Honest Prover Complexity

---

**Algorithm 4.12.1: Honest Prover Strategy**

**Input:** parameters $N \in \mathbb{N}$, $\tau \in (0, 0.01)$, black-box sample access to distribution $D$ over domain $[n]$, and tuple $(z_i)_{i \in [s]}$ such that for all $i \in [s]$, $z_i \in \mathrm{Supp}(D)$.

**Goal:** For every $i \in [s]$ such that $D(z_i) \geq \frac{\tau}{N}$ output $\pi(z_i) \in (0, 1)$ such that $\pi(z_i) \in D(z_i) \cdot [(1 - \tau/10, 1 + \tau/10)]$.

1. Draw $t = \frac{300 N \log N}{\tau^3}$ samples by $D$. Denote the sample as $(a_1, \ldots, a_t)$.

2. For every elements $x \in [N]$ denote $\widehat{D}(x) = \frac{|\{i \in [t]: a_i = x\}|}{t}$.

3. For every $i \in [s]$ set $\pi(z_i) = \max\left\{\widehat{D}(z_i), \frac{\tau}{2N}\right\}$.

4. Output $(\pi(z_i))_{i \in [s]}$

---

We show that the prover strategy outlined in Algorithm 4.12.1 is an efficient strategy for the honest prover in Protocol 4.1.1.

**Claim 4.13.** *Fix a sample $(z_i)_{i \in [s]}$. Algorithm 4.12.1 takes $\widetilde{O}(N \tau^{-3})$ samples and satisfies the following condition: with probability at least $1 - o(1)$, for every $i \in [s]$, such that $D(z_i) \geq \frac{\tau}{N}$ it holds that:*

$$\pi(z_i) \in D(z_i) \cdot [1 - \tau/10, 1 + 10\tau]$$

*Proof.* Fix $x \in [N]$. For every $i \in [t]$, define $\mathbb{1}_{a_i = x}$ to be the indicator of the event $a_i = x$. Therefore, $\mathbb{E}_{a_i \sim D}[\mathbb{1}_{a_i = x}] = D(x)$, and through the linearity of expectation, $\mathbb{E}_{(a_i)_{i \in [t]} \sim D^t}\left[\frac{1}{t} \sum_{i \in [t]} \mathbb{1}_{a_i = x}\right] = D(x)$. Assume $D(x) \geq \frac{\tau}{N}$. Through the multiplicative Chernoff bound:

$$\Pr_{(a_i)_{i \in [t]} \sim D^t}\left(\left|\frac{1}{t} \sum_{i \in [t]} \mathbb{1}_{a_i = x} - D(x)\right| > \frac{\tau}{10} D(x)\right) \leq 2\exp\left(-\frac{\tau^2 \cdot D(x)}{300} \cdot t\right) \leq 2\exp\left(-\log N\right) = \frac{2}{N}$$

Where the last inequality is due to the assumption that $D(x) \geq \frac{\tau}{N}$ and the choice of $t$. Since the sample $(z_i)_{i \in [s]}$ contains $O(s)$ elements, taking the union bound over all elements in $(z_i)_{i \in [s]}$:

$$\Pr_{(a_i)_{i \in [s]} \sim D^t}\left(\exists i \text{ s.t. } D(z_i) \geq \frac{\tau}{N} \text{ and } \pi(z_i) \notin D(z_i) \cdot \left[1 - \frac{\tau}{10}, 1 + \frac{\tau}{10}\right]\right) \leq s \cdot \frac{2}{N} = o(1)$$

Moreover, the sample complexity and runtime of the algorithm are both $O(t) = \widetilde{O}\left(N \tau^{-3}\right)$. $\square$

**Remark 4.14** (Completeness using approximate probabilities)**.** *For simplicity of presentaiton, the protocol's completeness analysis assumed that the tags provided by the honest prover specified each sample's correct bucket. In reality, however, the polynomial-time honest prover can only compute multiplicative approximations to the true probabilities. With high probability all the approximation will be quite good, but this still leaves the possibility that the honest prover might tag a sample as falling in an adjacent bucket (this can only happen to elements that are close to the edges of their true bucket). The completeness analysis of Section 4.1.1 already contains sufficient slack to allow for such errors: we elaborate below.*

*In more detail: we have shown that the prover obtains with high probability good approximations (up to a multiplicative factor of $\tau/10$) of the probability of all elements in $(z_i)_{i\in[s]}$ with probability above $\tau/N$. However, this does not imply that all the elements will be placed into their correct buckets. It is possible that an element $x \in B_\ell^D$ will satisfy $\pi(x) \in \left[\frac{e^{(\ell+1)\tau}}{N}, \frac{e^{(\ell+2)\tau}}{N}\right)$ or $\pi(x) \in \left[\frac{e^{(\ell-1)\tau}}{N}, \frac{e^{\ell\tau}}{N}\right)$. This will happen to elements that are close to the edges of each bucket.*

*Our protocol can withstand such errors while maintaining completeness: even if all elements close to the margins of each bucket are wrongly placed in adjacent buckets, the verifier will still pass all the tests, and we are still guaranteed that $\pi(z_i) \in D(z_i)[1 - \tau/10, 1 + \tau/10]$.*

*The reason behind this is that even though the result is phrased for a bucket partition as defined in Definition 3.3, in actuality, the verifier in Protocol 4.1.1 simply requires the following conditions from the bucket-partition:*

- *Every element in $\mathrm{Supp}(D)$ is in exactly one bucket (so that no element is accounted for twice by aggregating according to buckets).*

- *If $x, y \in [N]$ are in the same bucket, it holds $D(x) \approx D(y)$. Specifically, we consider a bucket partition for which $\frac{D(y)}{D(x)} \in [e^{-\tau}, e^\tau]$.*

*The second point is important, as it allows the verifier to conclude from the empirical mass of each bucket, namely $v_j = \widehat{q}_j$, an approximation of its size, by considering $Nv_je^{-j\tau}$. This quantity is close, with high probability to $Nq_\ell e^{-j\tau}$, which is a close approximation of $\left|B_j^D\right|$.*

*These approximations already come with considerable slack: the true size of bucket $j$ is anywhere in the interval $\left[\frac{q_j}{e^{(j+1)\tau}/N}, \frac{q_j}{e^{j\tau}/N}\right]$, and $v_j$ is anywhere in the interval $[q_j(1 - \tau), q_j(1 + \tau)]$. This slack is multiplicative, of magnitude $\theta(\tau)$, and is accounted for in the protocol, which in fact, allows even more slack than required.*

*We interpret the honest's prover's approximations to the true probabilities as a fuzzy bucket partition of the domain, $\left\{B_\ell^{D,fuzz}\right\}_\ell$, where elements are assigned to "fuzzy" buckets according to the empirical probabilities learned by the honest prover in Algorithm 4.12.1. Thus, each element in the support of $D$ lies in exactly one bucket, and (w.h.p. over the honest prover's samples) for each bucket $\ell$, if $x \in B_\ell^{D,fuzz}$ then $D(x) \in \left[\frac{e^{(\ell-1/10)\tau}}{N}, \frac{e^{(\ell+1/10)\tau}}{N}\right]$. This partition of the domain satisfies the two conditions given above, with a slightly higher slack in the bucket range (not considerably larger than before). The extra slack that this partition needs is already accounted for in the protocol's completeness analysis.*

### 4.3.3 Soundness of Protocol 4.1.1

Following the ideas outlined in the protocol overview in Section 4.1, the soundness condition claims that no matter what strategy a cheating prover might employ, if the verifier accepted, then with high probability, the prover commits to a lower bound for the size of each alleged bucket. I.e. the verifier is guaranteed that the approximation of the size each bucket $j$ obtained from the prover's answers, namely, $\frac{v_j}{e^{j\tau}/N}$, is bounded from below by roughly $\sum_\ell \frac{\widehat{q}_\ell}{e^{\ell\tau}/N} \cdot x_{\ell,j}$. Recall:

- For every $j$, $w_j$ is the alleged empirical mass of the $j$'th bucket under the distribution $U_{[N]}$, while $m_\ell$ is the *true* empirical mass of the $\ell$'th bucket under $U_{[N]}$. By definition of $\{y_{\ell,j}\}_{\ell,j}$, it holds that $w_j = \sum_\ell m_\ell y_{\ell,j}$.

24

- As in the previous section, for (almost) every $\ell$, it holds $m_\ell \approx \frac{|B_\ell^D|}{N} \approx \frac{N q_\ell e^{-\ell \tau}}{N} \approx \widehat{q}_\ell e^{-\ell \tau}$ (recall that $\widehat{q}_\ell$ is strongly concentrated around $q_\ell$).

And at the end of an accepting run, we are also guaranteed through Step (3b) that $w_j \approx v_j e^{-j\tau}$. Putting all of these together, we get that at the end of an accepting run, with high probability:

$$v_j e^{-j\tau} \approx w_j \approx \sum_\ell m_\ell y_{\ell,j} \approx \sum_\ell \widehat{q}_\ell e^{-\ell \tau} y_{\ell,j}$$

This is very close to what we actually wish to obtain. We need to replace $y_{\ell,j}$ with $x_{\ell,j}$, and as a consequence, we also replace the approximate equality with a lower bound on $v_j e^{-j\tau}$. The way we relate $x_{\ell,j}$ to $y_{\ell,j}$ is based on the unavoidable connection between the prover's mistags on samples in $S_0$ and samples in $S_1$: we show that a prover that wishes to mistag the samples drawn by $D$ must also mistag samples drawn by $U_{[N]}$ in a similar pattern (see Observation 4.1 for a more detailed intuition).

Thus, the majority of this section focuses on relating with high probability $x_{\ell,j}$ and $y_{\ell,j}$ for all bucket pairs $(\ell, j)$. We show that for *significant buckets* $x_{\ell,j} \approx y_{\ell,j}$. For *small buckets* we show instead that $\widehat{q}_\ell e^{-\ell \tau} x_{\ell,j} - m_\ell y_{\ell,j} = \widetilde{O}(1/s)$, i.e. $m_\ell y_{\ell,j}$ cannot be significantly larger than $\widehat{q}_\ell e^{-\ell \tau} x_{\ell,j}$ (this connection is what compels us to replace the approximate equality with a lower bound).

Instead of analyzing the prover's response to a sample drawn according to the process described in the protocol, we analyze the prover's answer with respect to a sample $(z_i)_{i \in [s]}$, and bits $(b_i)_{i \in [s]}$ that are distributed in the same way as the sample and bits drawn by the verifier, but were produced differently.

**Definition 4.15.** *Consider the joint distribution $(Z, B)$ defined as follows: $Z \sim \frac{1}{2}D + \frac{1}{2}U_{[N]}$, and:*

$$B\big|_{Z=z} = \begin{cases} 0, & w.p. \ \frac{D(z)}{D(z)+1/N} \\ 1, & w.p. \ \frac{1/N}{D(z)+1/N} \end{cases}$$

Let $(z_1, \ldots, z_s)$ and $(b_1, \ldots, b_s)$ respectively be the sample and the bits drawn by $V$ in Protocol 4.1.1. Define the random variable produced by collecting them together $S = ((z_1, b_1), \ldots, (z_s, b_s))$. For the same $s$ set as in Protocol 4.1.1, consider the the random variable $S' = ((z'_1, b'_1), \ldots, (z'_s, b'_s))$, where $i \in [s]$, $(z'_i, b'_i)$ is drawn i.i.d. according to distribution $(Z, B)$.

**Claim 4.16.** *Let $S$ and $S'$ be the random variables as above, then:*

$$\Delta_{SD}\left(S, S'\right) = 0$$

*Proof.* Note that by definition, for every sample $i$ drawn as in Protocol 4.1.1, $(z_i, b_i) = (x, 1)$ with probability $\frac{1}{2} \cdot \frac{1}{N}$ and $(z_i, b_i) = (x, 0)$ with probability $\frac{1}{2}D(x)$, as either bit is chosen with probability $\frac{1}{2}$, and then $z$ is sampled according to either $Q$ or $U_{[N]}$. Consider next the probability that $(z'_i, b'_i) = (x, 0)$. By definition, this is the product of the probability that distribution $\left(\frac{1}{2}D + \frac{1}{2}U_{[N]}\right)$ yielded $x$, and the probability that the bit chosen then was 0, given that $x$ was sampled. This is:

$$\left(\frac{1}{2}D(x) + \frac{1}{2}U_{[N]}(x)\right) \cdot \frac{D(x)}{D(x) + \frac{1}{N}} = \frac{1}{2}D(x)$$

Likewise, the probability that $(z'_i, b'_i) = (x, 1)$ is $\frac{1}{2} \cdot \frac{1}{N}$.

And so, we conclude that for every $x \in [N]$ and $b \in \{0, 1\}$, $(x, b)$ is as likely to have been produced through the process outlined in Protocol 4.1.1 as it is through distribution $(Z, B)$. $\square$

The important difference between $S$ and $S'$ is that the bits are set *after* the samples from the distributions are drawn. Since the prover is not given the bits $(b_i)_{i \in [s]}$, its view in Protocol 4.1.1 doesn't allow it to distinguish whether the sample it receives was drawn through the true process in Protocol 4.1.1, or each sample was drawn independently from $(Z, B)$. Therefore, for sake of the analysis in this section, we fix some cheating prover strategy $P^*$, and analyze the variables $\{x_{\ell,j}\}_{\ell,j}$ and $\{y_{\ell,j}\}_{\ell,j}$ where we think of $(z_i)_i$ and $(b_i)_i$ as drawn i.i.d. by $(Z, B)$. Whatever we conclude from this analysis then transfers to any prover strategy in the protocol. Concretely, consider the following mental experiment:

---

**Mental Experiment 4.16.1: Alternative Production of the View of Protocol 4.1.1**

1. V: draws $(z_1, \ldots, z_s)$ according to $\frac{1}{2}D + \frac{1}{2}U_{[N]}$. And sends it to P.

2. P: for every $i \in [s]$, set $\pi(z_i)$ as in Protocol 4.1.1. Send $(\pi(z_i))_{i \in [s]}$ to V.

3. V: for every $i \in [s]$, draw $b_i \sim B\big|_{Z=z_i}$.

---

For every $i \in [s]$, let $\mathrm{tag}(z_i) = \lfloor \log (N \cdot \pi(z_i)) / \tau \rfloor$. Every prover response in the mental experiment induces the following sets for every bucket index $\ell$ and for all $j$, that are well defined *before* $(b_i)_{i \in [s]}$ are set:

$$A_{\ell,j} = \left\{ i \in [s] : z_i \in B_\ell^D, \mathrm{tag}(z_i) = j \right\}$$

Claiming that $x_{\ell,j}$ and $y_{\ell,j}$ are close for some pair of bucket indices $(\ell, j)$ can be now phrased as a claim about the set $A_{\ell,j}$: we wish to show that after setting $(b_i)_i$ in Mental Experiment 4.16.1 roughly $\frac{\frac{e^{\ell\tau}}{N}}{\frac{e^{\ell\tau}}{N} + \frac{1}{N}}$ fraction of $A_{\ell,j}$ will be with its respective bit set to 0, and the remaining roughly $\frac{\frac{1}{N}}{\frac{e^{\ell\tau}}{N} + \frac{1}{N}}$ fraction of $A_{\ell,j}$ will be set with respective bit 1.

**Definition 4.17.** *For every $x \in [N]$ define the likelihood ratio of $x$ to be $r_x = \frac{D(x)}{D(x)+1/N}$*

We show that any sufficiently large enough collection of samples $R \subseteq [s]$ drawn according to $\frac{1}{2}Q + \frac{1}{2}U_{[N]}$ will be partitioned to two sets $S_0 \cap R$ and $S_1 \cap R$, of sizes roughly $\sum_{i \in R} r_{z_i}$ and $\sum_{i \in R} (1 - r_{z_i})$ respectively, with high probability.

**Claim 4.18.** *Fix $(z_1, \ldots, z_s)$. Let $R \subseteq [s]$. Assume $\sum_{i \in R} r_{z_i} > \frac{1}{\tau^2}$. With probability of at least $1 - 2e^{-\frac{1}{3}\tau^2 \sum_{i \in R} r_{z_i}}$ over the choice of $(b_i)_{i \in R}$:*

$$\left| |R \cap S_0| - \sum_{i \in R} r_{z_i} \right| < \tau \sum_{i \in R} r_{z_i}$$

*Similarly, assuming $\sum_{i \in R} (1 - r_{z_i}) > \frac{1}{\tau^2}$, with probability at least $1 - 2e^{-\frac{1}{3}\tau^2 \sum_{i \in R}(1-r_{z_i})}$*

$$\left| |R \cap S_1| - \sum_{i \in R}(1 - r_{z_i}) \right| < \tau \sum_{i \in R}(1 - r_{z_i})$$

26

*Proof.* Define $\mathbb{1}_{b_i=0}$ to be the indicator that $b_i = 0$. Recall that once $(z_1, \ldots, z_s)$ is set, $\mathbb{E}[\mathbb{1}_{b_i=0}] = \Pr(B = 0 \mid Z = z_i) = \frac{D(z_i)}{D(z_i)+1/N} = r_{z_i}$. Therefore, by linearity of expectation:

$$\mathbb{E}\left[|R \cap S_0|\right] = \sum_{i \in R} \mathbb{E}\left[\mathbb{1}_{b_i=0}\right] = \sum_{i \in R} r_{z_i}$$

And so, to wrap up the claim we just need to show that the random variable $|R \cap S_0|$ is concentrated around its mean, and indeed, by the multiplicative Chernoff Bound, since every bit $b_i$ is chosen independently of the other:

$$\Pr\left(\left||R \cap S_0| - \sum_{i \in R} r_{z_i}\right| > \tau \sum_{i \in R} r_{z_i}\right) \leq 2e^{-\frac{1}{3}\tau^2 \sum_{i \in R} r_{z_i}}$$

Note that $|R \cap S_0|$ assumes integer values, and the above inequality might be wrong if $|R \cap S_0| = 0$. However, assuming $\sum_{i \in R} r_{z_i} > \frac{1}{\tau^2}$, with high probability, $|R \cap S_0| \neq 0$ and the inequality above holds.

An analogous argument applies for $|R \cap S_1|$. $\qquad\square$

We apply this claim over the sets $A_{\ell,j}$ and $S^\ell$, defined as such:

**Definition 4.19.** *For every $\ell$, define $S^\ell = \left\{i \in [s] : z_i \in B_\ell^D\right\}$.*

Recall that our first goal was to prove that $x_{\ell,j} \approx y_{\ell,j}$ under some condition on the indices. Note that by definition $\frac{|A_{\ell,j} \cap S_0|}{|S^\ell \cap S_0|} = x_{\ell,j}$, and likewise $\frac{|A_{\ell,j} \cap S_1|}{|S^\ell \cap S_1|} = y_{\ell,j}$. And so, we apply the previous claim over these sets:

**Claim 4.20.** *Fix $(z_1, \ldots, z_s)$, as well as some prover response $(tag(z_1), \ldots, tag(z_s))$, that induces sets $A_{\ell,j} \subseteq [s]$. Let $(\ell,j)$ be such that $\sum_{i \in A_{\ell,j}} r_{z_i}, \sum_{i \in A_{\ell,j}} (1 - r_{z_i}) > \frac{\log(1000b(N,\tau))}{\tau^2}$. For every such pair $(\ell,j)$, with probability at least $1 - \delta_{\ell,j}$ over the choice of $(b_i)_{i \in [s]}$, it holds that:*

$$\frac{|A_{\ell,j} \cap S_0|}{|S^\ell \cap S_0|} \in \frac{|A_{\ell,j}|}{|S^\ell|} \cdot \left[e^{-4\tau}, e^{4\tau}\right]$$

*As well as:*

$$\frac{|A_{\ell,j} \cap S_1|}{|S^\ell \cap S_1|} \in \frac{|A_{\ell,j}|}{|S^\ell|} \cdot \left[e^{-4\tau}, e^{4\tau}\right]$$

*Where $\delta_{\ell,j} = 4\left(e^{-\frac{1}{3}\tau^2 \sum_{i \in A_{\ell,j}} r_{z_i}} + e^{-\frac{1}{3}\tau^2 \sum_{i \in A_{\ell,j}} (1 - r_{z_i})}\right)$.*

*Proof.* Let $\ell$ be some bucket index as assumed in the claim above. By Claim 4.18, with probability at most $2e^{-\frac{1}{3}\tau^2 \sum_{i \in A_{\ell,j}} r_{z_i}}$:

$$\left||A_{\ell,j} \cap S_0| - \sum_{i \in A_{\ell,j}} r_{z_i}\right| \geq \tau \sum_{i \in A_{\ell,j}} r_{z_i} \tag{14}$$

27

Also with probability at most $2e^{-\frac{1}{3}\tau^2 \sum_{i \in S^\ell} r_{z_i}}$:

$$\left| \left| S^\ell \cap S_0 \right| - \sum_{i \in S^\ell} r_{z_i} \right| \geq \tau \sum_{i \in S^\ell} r_{z_i} \tag{15}$$

Since by definition $\left| S^\ell \right| \geq |A_{\ell,j}|$, as $A_{\ell,j} \subseteq S^\ell$, taking union bound on these two events, we learn that with probability at most $4e^{-\frac{1}{3}\tau^2 \sum_{i \in A_{\ell,j}} r_{z_i}}$ both conditions apply.

Also note that for every $i$ such that $z_i \in B_\ell^D$, $D(z_i) \in \left[ \frac{e^{\ell\tau}}{N}, e^\tau \cdot \frac{e^{\ell\tau}}{N} \right)$, which implies $r_{z_i} \in \left[ \frac{e^{\ell\tau}/N}{e^{\ell\tau}/N+1/N}, e^\tau \cdot \frac{e^{\ell\tau}/N}{e^{\ell\tau}/N+1/N} \right)$. And so:

$$\left| \sum_{i \in A_{\ell,j}} r_{z_i} - |A_{\ell,j}| \frac{e^{\ell\tau}/N}{e^{\ell\tau}/N+1/N} \right| \leq (1 - e^\tau) \cdot |A_{\ell,j}| \frac{e^{\ell\tau}/N}{e^{\ell\tau}/N+1/N} \tag{16}$$

As well as:

$$\left| \sum_{i \in S^\ell} r_{z_i} - \left| S^\ell \right| \frac{e^{\ell\tau}/N}{e^{\ell\tau}/N+1/N} \right| \leq (1 - e^\tau) \cdot \left| S^\ell \right| \frac{e^{\ell\tau}/N}{e^{\ell\tau}/N+1/N} \tag{17}$$

Putting Inequalities (14),(15),(16), and (17) together, we get that with probability at least $4e^{-\frac{1}{3}\tau^2 \sum_{i \in A_{\ell,j}} r_{z_i}}$ over choice of $(b_i)_{i \in [s]}$, the following inequalities hold:

$$|A_{\ell,j} \cap S_0| \leq (1 + \tau) \sum_{i \in A_{\ell,j}} r_{z_i} \leq e^\tau \cdot e^\tau |A_{\ell,j}| \cdot \frac{e^{\ell\tau}/N}{e^{\ell\tau}/N+1/K} = e^{2\tau} |A_{\ell,j}| \cdot \frac{e^{\ell\tau}/N}{e^{\ell\tau}/N+1/K}$$

$$|A_{\ell,j} \cap S_0| \geq (1 - \tau) \sum_{i \in A_{\ell,j}} r_{z_i} \geq e^{-\tau} \cdot e^{-\tau} |A_{\ell,j}| \cdot \frac{e^{\ell\tau}/N}{e^{\ell\tau}/N+1/K} = e^{-2\tau} |A_{\ell,j}| \cdot \frac{e^{\ell\tau}/N}{e^{\ell\tau}/N+1/K}$$

As well as:

$$\left| S^\ell \cap S_0 \right| \in \left| S^\ell \right| \frac{e^{\ell\tau}/N}{e^{\ell\tau}/N+1/K} \cdot \left[ e^{-2\tau}, e^{2\tau} \right]$$

We thus get that with probability at least $4e^{-\frac{1}{3}\tau^2 \sum_{i \in A_{\ell,j}} r_{z_i}}$, $\left| S^\ell \cap S_0 \right| \neq 0$ and:

$$\frac{|A_{\ell,j} \cap S_0|}{|S^\ell \cap S_0|} \in \frac{|A_{\ell,j}|}{|S^\ell|} \cdot \left[ e^{-4\tau}, e^{4\tau} \right]$$

An analogous argument applies for the intersections with $S_1$. Taking a union bound over both conditions yields the desired result. $\qquad\square$

In particular, note that for pairs of indices $(j, \ell)$ such that $A_{\ell,j}$ is large enough as to have in expectation sufficiently many samples assigned bit 0 and 1 (concretely, at least $\widetilde{\Omega}(1/\tau^2)$ samples in expectation for each category), it holds that $x_{\ell,j} \approx y_{\ell,j}$:

28

**Corollary 4.21.** *Every pair $(\ell, j)$ such that $\sum_{i \in A_{\ell,j}} r_{z_i}, \sum_{i \in A_{\ell,j}} (1 - r_{z_i}) > \frac{\log(1000b(N,\tau))}{\tau^2}$, with probability at least $1 - \delta_{\ell,j}$ over the choice of $(b_i)_{i \in [s]}$, where $\delta_{\ell,j}$ is defined as in the previous claim, it holds that:*

$$\frac{x_{\ell,j}}{y_{\ell,j}} \in [e^{-8\tau}, e^{8\tau}]$$

*Proof.* Let $\ell$ and $j$ be bucket indices satisfying the conditions of the corollary. By definition, $x_{\ell,j} = \frac{|A_{\ell,j} \cap S_0|}{|S^\ell \cap S_0|}$, as well as $y_{\ell,j} = \frac{|A_{\ell,j} \cap S_1|}{|S^\ell \cap S_1|}$. Therefore, Claim 4.20, we get $\frac{x_{\ell,j}}{y_{\ell,j}} \in [e^{-8\tau}, e^{8\tau}]$ □

Recall that the verifier checks that $w_j \approx v_j e^{-j\tau}$. By definition, $w_j = \sum_\ell m_\ell y_{\ell,j}$. And so, if the verifier accepted, then:

$$v_j e^{-j\tau} \approx \sum_\ell m_\ell y_{\ell,j}$$

In order to prove that the soundness condition holds, we need to show that for every $j$, $\sum_\ell m_\ell y_{\ell,j}$ is larger than roughly $\sum_\ell \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau}$. In order to show this, we show that with high probability, for all $j$ and every $\ell$, one of the following must hold:

- $m_\ell y_{\ell,j} \approx q_\ell e^{-\ell\tau} x_{\ell,j}$

- $m_\ell y_{\ell,j}$ is not significantly smaller than $q_\ell e^{-\ell\tau} x_{\ell,j}$.

Starting with the first condition, the pairs of indices $(\ell, j)$ for which the former holds are those for which $x_{\ell,j} \approx y_{\ell,j}$. Formally:

**Claim 4.22.** *With probability $0.99$ over the choice of $(z_1, \ldots, z_s)$, any prover response $(\pi(z_1), \cdots \pi(z_s))$, induces sets $\{A_{\ell,j}\}_{\ell,j}$, and satisfies the condition that with probability at least $0.98$ over the choice of $(b_i)_{i \in [s]}$, for all bucket indices pairs $(\ell, j)$ that satisfy $\sum_{i \in A_{\ell,j}} r_{z_i}, \sum_{i \in A_{\ell,j}} (1 - r_{z_i}) > \frac{\log(1000b(N,\tau))}{\tau^2}$, it holds that:*

$$\left| m_\ell \cdot y_{\ell,j} - \widehat{q}_\ell e^{-\ell\tau} x_{\ell,j} \right| < 15\tau \cdot \widehat{q}_\ell e^{-\ell\tau} x_{\ell,j}$$

*Proof.* For all bucket indices pairs $(\ell, j)$ that satisfy the conditions of the claim, the following apply:

- We conclude that $\left| B_\ell^D \right|$ has to be big: since $A_{\ell,j} \subseteq S^\ell$, the fact that $\sum_{i \in A_{\ell,j}} (1 - r_{z_i}) > \frac{\log(1000b(N,\tau))}{\tau^2}$, implies that $\sum_{i \in S^\ell} (1 - r_{z_i}) > \frac{\log(1000b(N,\tau))}{\tau^2}$, and since $\mathbb{E}[s_1 \cdot m_\ell] = \sum_{i \in S^\ell} (1 - r_{z_i}) = \frac{|B_\ell^D|}{N}$. With overwhelming probability $s_1 \in [s/3, 2s/3]$ (see Claim 4.3.1), and so, with high probability:

$$\frac{\left| B_\ell^D \right|}{N} \geq \frac{3}{2s} \sum_{i \in S^\ell} (1 - r_{z_i}) \geq \frac{\log(1000b(N,\tau))}{s\tau^2}$$

By Claim 4.10, we conclude that with probability at least $0.99$, every $\ell$ that satisfies the above conditions also satisfies:

$$\left| m_\ell - q_\ell e^{-\ell\tau} \right| \leq 2\tau q_\ell e^{-\ell\tau}$$

And so:

$$\left| m_\ell x_{\ell,j} - q_\ell e^{-\ell\tau} x_{\ell,j} \right| < 2\tau \cdot q_\ell e^{-\ell\tau} x_{\ell,j} \tag{18}$$

- By Corollary 4.21, with probability at least 0.99, over all buckets satisfying:

$$\sum_{i \in A_{\ell,j}} r_{z_i}, \sum_{i \in A_{\ell,j}} (1 - r_{z_i}) > \frac{\log\left(1000 b(N, \tau)\right)}{\tau^2}$$

We get that $\frac{y_{\ell,j}}{x_{\ell,j}} \in [e^{-8\tau}, e^{8\tau}]$. This implies:

$$|m_\ell y_{\ell,j} - m_\ell x_{\ell,j}| < \left(e^{8\tau} - 1\right) m_\ell x_{\ell,j} < 9\tau \cdot m_\ell x_{\ell,j} \tag{19}$$

Therefore, taking the union bound, with probability at least 0.98, for all $\ell$ and $j$ satisfying the conditions in the statement, by putting Inequalities (18) and (19) together, we conclude that:

$$\left|m_\ell y_{\ell,j} - q_\ell e^{-\ell\tau} x_{\ell,j}\right| \leq |m_\ell y_{\ell,j} - m_\ell x_{\ell,j}| + \left|m_\ell x_{\ell,j} - q_\ell e^{-\ell\tau} x_{\ell,j}\right| \leq 9\tau m_\ell x_{\ell,j} + 2\tau q_\ell e^{-\ell\tau} x_{\ell,j} \leq 12\tau q_\ell e^{-\ell\tau} x_{\ell,j}$$

Where the last inequality is justified through Inequality (18). Finally, observing that $\sum_{i \in A_{\ell,j}} r_{z_i} > \frac{\log(1000 b(N,\tau))}{\tau^2}$, also implies as above that $q_\ell > \frac{\log(1000 b(N,\tau))}{s\tau^2}$. We get through Claim 4.9 that $|q_\ell - \widehat{q}_\ell| \leq (e^\tau - 1) q_\ell$. Plugging this in the inequality above yields the desired result. $\qquad\square$

Moving on to pairs $(\ell, j)$ that don't satisfy $\sum_{i \in A_{\ell,j}} r_{z_i}, \sum_{i \in A_{\ell,j}} (1 - r_{z_i}) > \frac{\log(1000 b(N,\tau))}{\tau^2}$, we relate the quantities $m_\ell y_{\ell,j}$ and $q_\ell e^{-\ell\tau} x_{\ell,j}$ additively, using Markov's inequality.

**Claim 4.23.** *Fix* $(z_1, \ldots, z_s)$, *as well as prover's response* $(\pi(z_1), \ldots, \pi(z_s))$, *which induces sets* $\{A_{\ell,j}\}_{\ell,j}$. *With probability at least* 0.99, *every pair* $(\ell, j)$ *for which* $\sum_{i \in A_{\ell,j}} (1 - r_{z_i}) \leq \frac{\log(1000 b(N,\tau))}{\tau^2}$ *satisfies:*

$$m_\ell y_{\ell,j} - \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau} \geq -\frac{1}{s} \cdot \frac{200 b(N,\tau)^2 \cdot \log\left(1000 b(N,\tau)\right)}{\tau^2} \tag{20}$$

*Proof.* Since $m_\ell y_{\ell,j} \geq 0$, suffice to show that with probability of at least 0.99 over the choice of $b \in \{0,1\}^s$ for all $(\ell, j)$ such that $\sum_{i \in A_{\ell,j}} (1 - r_{z_i}) \leq \frac{\log(1000 b(N,\tau))}{\tau^2}$, it holds that:

$$\widehat{q}_\ell x_{\ell,j} e^{-\ell\tau} \leq \frac{1}{s} \cdot \frac{200 b(N,\tau)^2 \cdot \log\left(1000 b(N,\tau)\right)}{\tau^2}$$

Indeed, by Definition 4.17, for every $i \in B_\ell^D$, $r_{z_i} \leq e^{(\ell+1)\tau} (1 - r_{z_i})$, and by extension:

$$\mathbb{E}_{b_i \sim B}\Big|_{Z=z_i} [s_0 \widehat{q}_\ell x_{\ell,j}] = \sum_{i \in A_{\ell,j}} r_{z_i} \leq e^{(\ell+1)\tau} \sum_{i \in A_{\ell,j}} (1 - r_{z_i}) \leq e^{(\ell+1)\tau} \frac{\log\left(1000 b(N,\tau)\right)}{\tau^2}$$

From which we get that by Markov's Inequality, with probability at most $\frac{1}{100 b(N,\tau)^2}$, it holds that:

$$s_0 \widehat{q}_\ell x_{\ell,j} \geq 100 b(N,\tau)^2 \cdot e^{(\ell+1)\tau} \frac{\log\left(1000 b(N,\tau)\right)}{\tau^2} \tag{21}$$

Taking union bound over all such buckets, with probability at least 0.99, all pairs $(\ell, j)$ for which $\sum_{i \in A_{\ell,j}} (1 - r_{z_i}) \leq \frac{\log(1000 b(N,\tau))}{\tau^2}$, also satisfy:

$$s_0 \widehat{q}_\ell x_{\ell,j} \leq 100 b(N,\tau)^2 \cdot e^{(\ell+1)\tau} \frac{\log\left(1000 b(N,\tau)\right)}{\tau^2}$$

30

Since with high probability $s_0 \in \frac{s}{2} \cdot [e^{-\tau}, e^{\tau}]$, we conclude that for all such pairs, with probability at least $0.99$:

$$\widehat{q}_\ell e^{-\ell\tau} x_{\ell,j} \leq \frac{100 e^\tau b(N,\tau)^2}{s_0} \cdot \frac{\log\left(1000 b(N,\tau)\right)}{\tau^2} \leq \frac{1}{s} \cdot \frac{200 b(N,\tau)^2 \cdot \log\left(1000 b(N,\tau)\right)}{\tau^2}$$

$\square$

In order to account for all possible indices pairs, we are left to deal with those indices $(\ell, j)$ for which $\sum_{i \in A_{\ell,j}} r_{z_i} \leq \frac{\log(1000 b(N,\tau))}{\tau^2}$. We focus our attention only to the case where $\frac{e^{\ell\tau}}{N} > \frac{\tau^2}{N}$.

**Claim 4.24.** *Fix $(z_i)_{i \in [s]}$, as well as prover's response $(\pi(z_i))_{i \in [s]}$, which induces sets $\{A_{\ell,j}\}_{\ell,j}$. W.p. at least $0.99$ over all $(\ell, j)$ such that $\sum_{i \in A_{\ell,j}} r_{z_i} \leq \frac{\log(1000 b(N,\tau))}{\tau^2}$ and $\frac{e^{\ell\tau}}{N} \geq \frac{\tau}{N}$, it holds that:*

$$m_\ell y_{\ell,j} - \widehat{q} e^{-\ell\tau} x_{\ell,j} \geq -\frac{1}{s} \cdot \frac{200 b(N,\tau)^2 \cdot \log\left(1000 b(N,\tau)\right)}{\tau^4}$$

*Proof.* This proof follows the same line of reasoning as the proof of Claim 4.23. First, observe that:

$$\mathbb{E}_{b_i \sim B}\Big|_{Z = z_i} \left[ s_0 \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau} \right] = e^{-\ell\tau} \sum_{i \in A_{\ell,j}} r_{z_i} \leq \frac{1}{\tau^2} \cdot \frac{\log\left(1000 b(N,\tau)\right)}{\tau^2} \leq \frac{\log\left(1000 b(N,\tau)\right)}{\tau^4} \tag{22}$$

Next, as in Claim 4.23, applying both Markov's inequality with the union bound, alongside the fact that $s_0 \in \frac{s}{2}[e^{-\tau}, e^\tau]$ with high probability, yields the desired result. $\square$

**Claim 4.25** (Soundness of Protocol 4.1.1). *No matter what cheating strategy a cheating prover might employ, with probability at least $0.95$ over the samples of the verifier, either the verifier rejects, or for every $j$ such that $v_j \geq \frac{\tau^2}{\log N}$:*

$$v_j e^{-j\tau} \geq (1 - 18\tau) \sum_{\ell : e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell e^{-\ell\tau} x_{\ell,j} - \mathsf{poly}\left(\log N, \tau^{-1}\right) \cdot \frac{1}{s}$$

*Proof.* Assume a run of Protocol 4.1.1 terminated with the verifier accepting. Let $(z_1, \ldots, z_s)$ be the samples drawn by the verifier and sent to the prover, and $(\pi(z_1), \ldots \pi(z_s))$ the prover's response that induces sets $A_{\ell,j}$. Since the verifier accepted, we know that for every $j$ such that $v_j \geq \frac{\tau^2}{\log N}$:

$$\left| w_j - v_j e^{-j\tau} \right| \leq 4\tau v_j e^{-j\tau} \leq 5\tau w_j \tag{23}$$

Denote $\mathtt{GOOD} = \left\{ (\ell, j) : \sum_{i \in A_{\ell,j}} r_{z_i}, \sum_{i \in A_{\ell,j}} (1 - r_{z_i}) \leq \frac{\log(1000 b(N,\tau))}{\tau^2}, \right\}$, and $\mathtt{BAD}$ the collection of all other pairs. By definition, $w_j = \sum_\ell m_\ell y_{\ell,j}$. By Claims 4.22, 4.23, and 4.24, with probability at least $0.95$ over $(z_1, \ldots z_s)$, the prover's randomness, and $b \in \{0,1\}^s$, we get:

$$w_j = \sum_\ell m_\ell y_{\ell,j}$$

$$\geq \sum_{\ell:(\ell,j) \in \mathtt{GOOD}} m_\ell y_{\ell,j} + \sum_{\ell:(\ell,j) \in \mathtt{BAD}, e^{\ell\tau} \geq \tau^2} m_\ell y_{\ell,j}$$

$$\geq (1 - 15\tau) \sum_{\ell:(\ell,j) \in \mathtt{GOOD}} \widehat{q}_\ell e^{-\ell\tau} x_{\ell,j} + \sum_{\ell:(\ell,j) \in \mathtt{BAD}, e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell e^{-\ell\tau} x_{\ell,j} - \mathsf{poly}(\log N, \tau^{-1}) \cdot \frac{1}{s}$$

$$\geq (1 - 15\tau) \sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell e^{-\ell\tau} x_{\ell,j} - \mathsf{poly}(\log N, \tau^{-1}) \cdot \frac{1}{s}$$

$\square$

# 5 Collisions Matching Test

In the previous section, a verifier with black-box sample access to a distribution $D$ over domain $[N]$ drew $s = \widetilde{O}\left(\sqrt{N}\right)$ poly $\left(\tau^{-1}\right)$ samples $(z_i)_{i \in [s]}$ for some $\tau \in (0, 0.01)$, and obtained the colletion of claims $(\pi(z_i))_{i \in [s]}$, which allegedly approximates $(D(z_i))_{i \in [s]}$. In this section, we present an algorithm (tester) that given sample access to $D$, the input $((z_i, \pi(z_i)))_{i \in [s]}$, as well as parameters $N \in \mathbb{N}$ and $\sigma \in (0, 1)$, where $\sigma$ is assumed to satisfy $\sigma = \Omega(\sqrt{\tau})$, satisfies the following conditions:

**Tester Completeness.** If for all $i \in [s]$, $\pi(z_i) \approx D(z_i)$ the tester will accept with high probability

**Tester Soundness.** If the claims $(\pi(z_i))_{i \in [s]}$ are $\sigma$-far from $(D(z_i))_{i \in [s]}$ in the following sense: $\frac{1}{s} \sum_{i \in [s]} \left(1 - \min\left\{\frac{D(z_i)}{\pi(z_i)}, \frac{\pi(z_i)}{D(z_i)}\right\}\right) \geq \sigma$, then with high probability either the tester rejects, or the following condition holds: let $\{v_j\}_j$, $\{\widehat{q}_\ell\}_\ell$ and $\{x_{\ell,j}\}_{\ell,j}$ induced by $((z_i, \pi(z_i)))_{i \in [s]}$ be defined as explained in Section 4. There must be some bucket index $j$ with significant alleged mass for which $Nv_j e^{-j\tau}$ (i.e. the alleged size of the bucket), satisfies:

$$Nv_j e^{-j\tau} \leq \left(1 - \Omega(\sigma^2)\right) \sum_{\ell:e^{\ell\tau} \geq \tau^2} N\widehat{q}_\ell x_{\ell,j} e^{-\ell\tau} \tag{24}$$

Recall that in the previous section we showed that if the run of the *Bucket Size Verification Protocol* terminated with the verifier accepting, it must be that for all $j$ with significant alleged mass, the alleged size of the $j$'th $(N, \tau)$-bucket, $Nv_j e^{-j\tau}$, is *lower bounded* by $(1-O(\tau)) \sum_{\ell:e^{\ell\tau} \geq \tau^2} N\widehat{q}_\ell x_{\ell,j} e^{-\ell\tau}$. And so, conflating these two conditions one against the other later in Section 6, we show that it is possible to choose $\tau = O(\sigma^2)$ and $s = \widetilde{O}(\sqrt{N})$poly$(\sigma^{-1})$, so that if the prover provided tags that are $\sigma$-*far* from the truth in the above-mentioned sense, it cannot be that with high probability both the *Bucket Size Verification Protocol* and the tester presented in this section pass with high probability, as if they had, there would be some bucket $j$ for which $Nv_j e^{-j\tau}$ will be bounded from above by $\left(1 - \Omega(\sigma^2)\right) \sum_{\ell:e^{\ell\tau} \geq \tau^2} N\widehat{q}_\ell x_{\ell,j} e^{-\ell\tau}$, which is smaller than its lower bound of $(1 - O(\tau)) \sum_{\ell:e^{\ell\tau} \geq \tau^2} N\widehat{q}_\ell x_{\ell,j} e^{-\ell\tau}$, reaching contradiction. More on the combination of the protocol and the tester in Section 6.

**Tester Outline.** Tester 5.1.1 works as follows: it produces from the tagged sample $(z_i, \pi(z_i))$ the histogram $\{v_j\}_j$ as explained in Protocol 4.1.1. If the tags are correct, we expect that for all $\ell$ $v_j = \widehat{q}_\ell \approx q_\ell = D(B_\ell^D)$. Then, it draws $s$ fresh samples from $D$, $T = (t_1, \ldots, t_s)$. By choice of $s$, there will be a lot of collisions between elements of $(z_i)_{i \in [s]}$ and $(t_i)_{i \in [s]}$. For every alleged bucket $j$, we define the following variables for counting the collisions involving samples tagged as belonging to bucekt $j$:

$$\widetilde{C}_j = |\{(k, m) \in [s] \times [s] : z_k = t_m, \text{tag}(z_k) = j\}|$$

To estimate the expected value of this variable, consider that there are $s \cdot v_j$ samples in $S$ tagged as belonging to bucket $j$. If the tags are correct, each sample is of probability approximately $e^{j\tau}/N$, and so, we expect that:

$$\mathbb{E}_{T \sim D^s}\left[\widetilde{C}_j\right] = (s \cdot v_j) \cdot s \cdot \frac{e^{\ell\tau}}{N} = \frac{s^2}{N} v_j e^{\ell\tau}$$

And indeed, this is exactly what the tester checks, it goes over every (significantly heavy) bucket $j$, and checks that $\widetilde{C}_j \approx \frac{s^2}{N} v_j e^{j\tau}$. If the tags are correct, from concentration bounds we get that this check passes with high probability for all $j$. If the tags are far from being correct, i.e. $\frac{1}{s} \sum_{i \in [s]} \left( 1 - \min \left\{ \frac{D(z_i)}{\pi(z_i)}, \frac{\pi(z_i)}{D(z_i)} \right\} \right) \geq \sigma$, then, we argue that it must be that there exists some (significantly heavy) bucket $j$ for which: $N v_j e^{-j\tau} \leq \left( 1 - \Omega(\sigma^2) \right) \sum_{\ell : e^{\ell\tau} \geq \tau^2} N \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau}$. This is because the collision matching test rejects with high probability any tags that don't satisfy the following approximate equality for every $j$:

$$v_j \frac{e^{j\tau}}{N} \approx \sum_\ell \widehat{q}_\ell x_{\ell,j} \frac{e^{\ell\tau}}{N}$$

And we show that satisfying this condition on every bucket, implies an upper bound on the alleged size of some bucket $j$ (this is inspired by techniques in [HR22], and illustrated in Section 2).

**Proposition 5.1.** *There exists a tester $\mathcal{T}$ that gets as input parameters $\sigma \in (0, 0.1), N \in \mathbb{N}$, sample access to a distribution $D$ over domain $[N]$, and a sample $(z_i)_{i \in [s]}$ that was drawn i.i.d. by $D$ of size $s = \widetilde{\Theta} \left( \sqrt{N} \right) \mathsf{poly}(1/\tau)$, where $\tau = O(\sigma^2)$, alongside the tuple $(\pi(z_i))_{i \in [s]}$, where $\pi(z_i) \in (0, 1]$ for all $i \in [s]$. $D$ is assumed to satisfy $D(x) \leq \frac{\tau}{\sqrt{N}}$ for all $x \in [N]$. The tester's sample complexity and runtime are both $\widetilde{O} \left( \sqrt{N} \right) \mathsf{poly}(\tau^{-1})$, and at the end of the run, the following apply:*

- ***Completeness.*** *If for all $i \in [s]$, $\frac{\pi(z_i)}{D(z_i)} \in [e^{-\tau}, e^{\tau}]$, then with probability at least $0.9$ over the choice of $(z_i)_{i \in [s]}$, and the samples drawn by $\mathcal{T}$, $\mathcal{T}$ accepts.*

- ***Soundness.*** *Let $\{v_j\}_j$ be the $(N, \tau)$-histogram induced by $((z_i, \pi(z_i)))_{i \in [s]}$ (see Protocol 4.1.1 for details). If:*

$$\frac{1}{s} \sum_{i \in [s]} \left( 1 - \min \left\{ \frac{\pi(z_i)}{D(z_i)}, \frac{D(z_i)}{\pi(z_i)} \right\} \right) > \sigma$$

*Then, with probability at least $0.9$ over $(z_i)_{i \in [s]}$ and $\mathcal{T}$'s samples, there exists at least one bucket index $j$ such that $\frac{e^{j\tau}}{N} \geq \frac{\tau}{N}$, $v_j \geq \frac{\tau^2}{\log N}$ and:*

$$\sum_{\ell : e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau} \geq \left( 1 + 0.2\sigma^2 - 10\tau \right) v_j e^{-j\tau} \tag{25}$$

*Where $\{v_j\}_j$ and $\{x_{\ell,j}\}_{\ell,j}$ are induced by $((z_i, \pi(z_i)))_{i \in [s]}$ as defined in Definitions 4.7 and 4.2 respectively.*

We show that Tester 5.1.1 meets the conditions specified in Proposition 5.1.

## 5.1 Completeness of Tester 5.1.1

The proof that Tester 5.1.1 satisfies the completeness condition outlined in Proposition 5.1 follows Herman and Rothblum [HR22], and is restated in Appendix A, Claim A.2.

We thus proceed to show that Tester 5.1.1 satisfies the soundness condition in Proposition 5.1.

---

**Tester 5.1.1: Collisions Matching Tester**

**Input:** sample access to distribution $D$ over domain $[N]$, parameters $N \in \mathbb{N}$, $\sigma, \tau \in (0,1)$ such that $\tau = O(\sigma^2)$, as well as a tagged sample $S = ((z_i, \pi(z_i)))_{i \in [s]}$, for $s = \widetilde{O}(\sqrt{N})\mathsf{poly}(\tau^{-1})$, such that $(z_i)_{i \in [s]}$ was drawn i.i.d. by $D$. Assume that for all $x \in [N]$ $D(x) \leq \frac{\tau}{\sqrt{N}}$.

1. Set $\{v_j\}_j$ and $(\mathrm{tag}(z_i))_{i \in [s]}$ as in Protocol 4.1.1.

2. Draw $s$ fresh samples from $D$, denote this sample as $T = (t_1, t_2, \ldots, t_s)$. For every $(N, \tau)$-bucket index $j$, define:
$$\widetilde{C}_j = |\{(k, m) \in [s] \times [s] : z_k = t_m, \mathrm{tag}(z_k) = j\}|$$

3. Reject unless for every bucket $j$ such that $v_j \geq \frac{\tau^2}{\log N}$, and $\frac{e^{j\tau}}{N} \geq \frac{\tau}{N}$, it holds that:
$$\left| \widetilde{C}_j - \frac{s^2}{N} v_j e^{j\tau} \right| \leq 5\tau \frac{s^2}{N} v_j e^{j\tau}$$

---

## 5.2 Soundness of Tester 5.1.1

Let $\{x_{\ell,j}\}_{\ell,j}$ be defined with respect to the tagged sample $((z_i, \pi(z_i)))_{i \in [s]}$ as in Definition 4.2.

We wish to show that if $\frac{1}{s} \sum_{i \in [s]} \left(1 - \min\left\{\frac{D(z_i)}{\pi(z_i)}, \frac{\pi(z_i)}{D(z_i)}\right\}\right)$ then there exists some bucket index $j$ such that $v_j \geq \frac{\tau^2}{\log N}$ and $e^{j\tau} \geq \tau$ that satisfies Inequality (25). We actually show a slightly stronger condition. We show that under the above mentioned assumption there exists with high probability some bucket index $j$ that satisfies $v_j \geq \frac{\tau^2}{\log N}$, $e^{j\tau} \geq \tau$, as well as $\widehat{q}_L x_{L,j} \leq \tau v_j$, for which Inequality (25) holds. Recall that by definition $\widehat{q}_L$ is the empirical mass of $B_L^D = \left\{x \in [N] : D(x) \leq \frac{\tau^2}{N}\right\}$. Since $q_L \leq N \cdot \frac{\tau^2}{N} \leq \tau^2$, we expect many alleged buckets to satisfy this condition. We define:

**Definition 5.2.** *Fix sample $(z_i)_{i \in [s]}$ as well as prover response $(\pi(z_i))_{i \in [s]}$, and let $\{v_j\}_j$ be the alleged $(N, \tau)$-histogram induced by the tagged sample. Denote the set of* good *bucket indices $J = \left\{j : v_j \geq \frac{\tau^2}{\log N}, e^{j\tau} \geq \tau, \widehat{q}_L x_{L,j} \leq \tau v_j\right\}$*

We first observe that by definition we can rewrite the soundness condition as follows:

**Claim 5.3.** *If $\frac{1}{s} \sum_{i \in [s]} \left(1 - \min\left\{\frac{\pi(z_i)}{D(z_i)}, \frac{D(z_i)}{\pi(z_i)}\right\}\right) > \sigma$, then, given that $\tau < \sigma^2$ and $\sigma < 0.1$, it holds that:*
$$\sum_j \sum_\ell \widehat{q}_\ell x_{\ell,j} \left(1 - e^{-|\ell - j|\tau}\right) \geq 0.99\sigma$$

*Proof.*

$$\frac{1}{s}\sum_{i\in[s]}\left(1-\min\left\{\frac{\pi(z_i)}{D(z_i)},\frac{D(z_i)}{\pi(z_i)}\right\}\right)=\frac{1}{s}\sum_{j}\sum_{i:\text{tag}(z_i)=j}\left(1-\min\left\{\frac{\pi(z_i)}{D(z_i)},\frac{D(z_i)}{\pi(z_i)}\right\}\right)$$

$$=\frac{1}{s}\sum_{j}\sum_{\ell}\sum_{i:\text{tag}(z_i)=j,z_i\in B_\ell^D}\left(1-\min\left\{\frac{\pi(z_i)}{D(z_i)},\frac{D(z_i)}{\pi(z_i)}\right\}\right)$$

$$\leq e^{2\tau}\frac{1}{s}\sum_{j}\sum_{\ell}s\widehat{q}_\ell x_{\ell,j}\left(1-\min\left\{\frac{e^{j\tau}/N}{e^{\ell\tau}/N},\frac{e^{\ell\tau}/N}{e^{j\tau}/N}\right\}\right)$$

$$\leq e^{2\tau}\sum_{j}\sum_{\ell}\widehat{q}_\ell x_{\ell,j}\left(1-\min\left\{\frac{e^{j\tau}/N}{e^{\ell\tau}/N},\frac{e^{\ell\tau}/N}{e^{j\tau}/N}\right\}\right)$$

$$=e^{2\tau}\sum_{j}\sum_{\ell}\widehat{q}_\ell x_{\ell,j}\left(1-e^{-|\ell-j|\tau}\right)$$

Where the first inequality is due to the fact that by definition for every $z\in B_\ell^D$ and $\text{tag}(z)=j$ it holds that both $D(z)\in\left[\frac{e^{\ell\tau}}{N},\frac{e^{(\ell+1)\tau}}{N}\right)$ and $\pi(z)\in\left[\frac{e^{j\tau}}{N},\frac{e^{(j+1)\tau}}{N}\right)$, as well as since the number of samples $i\in[s]$ satisfying $z_i\in B_\ell^D$ and $\text{tag}(z_i)=j$ is by definition $s\widehat{q}_\ell x_{\ell,j}$. We conclude that since $\sigma<0.1$ and $\tau<\sigma^2$ that the soundness conditions implies that:

$$\sum_{j}\sum_{\ell}\widehat{q}_\ell x_{\ell,j}\left(1-e^{-|\ell-j|\tau}\right)\geq 0.99\sigma$$

$\square$

From now on, we will consider the soundness condition to be $\sum_{j}\sum_{\ell}\widehat{q}_\ell x_{\ell,j}\left(1-e^{-|\ell-j|\tau}\right)\geq 0.99\sigma$.

Next, we show it must be that some $j\in J$ for which the sum $\sum_{\ell}\widehat{q}_\ell x_{\ell,j}\left(1-e^{-|\ell-j|\tau}\right)$ is relatively large, i.e. there are alleged buckets that satisfy the conditions of set $J$ and account for a significant portion of the distance:

**Claim 5.4.** *Let $\sigma\in(0,0.01)$ be such that $\tau\leq\sigma^2$. Assume $\sum_{i\in[s]}\left(1-\min\left\{\frac{\widetilde{Q}(z_i)}{Q(z_i)},\frac{Q(z_i)}{\widetilde{Q}(z_i}\right\}\right)>\sigma$. With probability at least $0.99$ over the choice of $(z_i)_{i\in[s]}$, there exists some $j\in J$ such that:*

$$\sum_{\ell}\widehat{q}_\ell x_{\ell,j}\left(1-e^{-|\ell-j|}\right)>0.8v_j\sigma$$

*Proof.* As explained above, assuming the soundness condition implies that:

$$\sum_{j}\sum_{\ell}\widehat{q}_\ell x_{\ell,j}\left(1-e^{-|\ell-j|}\right)\geq 0.99\sigma$$

Denote $J_1=\left\{j\ :\ v_j\leq\frac{\tau^2}{\log N}\right\}$. It follows immediately that $\sum_{j\in J_1}\sum_{\ell}\widehat{q}_\ell x_{\ell,j}=\sum_{j\in J_1}v_j\leq b(N,\tau)\cdot\frac{\tau^2}{\log N}\leq\tau$.

Similarly, we define $J_2 = \{j : \hat{q}_L x_{L,j} > \tau v_j\}$. Observe that for every $j \in J_2$:

$$\sum_{\ell:e^{\ell\tau} \geq \tau^2} \hat{q}_\ell x_{\ell,j} \leq v_j \leq \frac{\hat{q}_L x_{L,j}}{\tau}$$

Therefore, with probability at least 0.99 over the choice of $(z_i)_{i\in[s]}$:

$$\sum_{j\in J_2} \sum_\ell \hat{q}_\ell x_{\ell,j} \leq \sum_{j\in J_2} \left( \sum_{\ell:e^{\ell\tau} \geq \tau^2} \hat{q}_\ell x_{\ell,j} + \hat{q}_L x_{L,j} \right)$$

$$\leq \sum_{j\in J_2} (\hat{q}_L x_{L,j}/\tau + \hat{q}_L x_{L,j})$$

$$\leq \left( \frac{1}{\tau} + 1 \right) \sum_{j\in J_2} \hat{q}_L x_{L,j}$$

$$\leq \left( \frac{1}{\tau} + 1 \right) \hat{q}_L$$

$$\leq \left( \frac{1}{\tau} + 1 \right) 100 q_L$$

$$\leq \left( \frac{1}{\tau} + 1 \right) 100\tau^2$$

$$\leq 100\tau + 100\tau^2$$

$$\leq 0.1\sigma$$

Where the third to last inequality is due to Markov's Inequality, since $\mathbb{E}_{(z_i)_{i\in[s]}}[\hat{q}_L] = q_L$; the second to last from the fact that since the domain of $D$ is $[N]$, then $q_L \leq N \cdot \frac{\tau^2}{N} \leq \tau^2$, and the last inequality is due to the assumption over $\sigma$ and $\tau$. Finally, define $J_3 = \{j : e^{j\tau} < \tau\}$:

$$\sum_{j\in J_3} \sum_\ell \hat{q}_\ell x_{\ell,j} \leq \sum_{j\in J_3} \sum_\ell \hat{q}_\ell x_{\ell,j} \leq \sum_{j\in J_3} v_j \leq 2\tau \leq 0.01\sigma$$

Since $J = \{j : j \notin J_1 \cup J_2 \cup J_3\}$, we conclude:

$$\sum_{j\in J} v_j = \sum_{j\in J} \sum_\ell \hat{q}_\ell x_{\ell,j} \geq 1 - 0.15\sigma$$

Since for all $\ell$ and $j$, $\left(1 - e^{-|\ell-j|}\right) \leq 1$, we get:

$$0.99\sigma \leq \sum_j \sum_\ell \hat{q}_\ell x_{\ell,j} \left(1 - e^{-|\ell-j|}\right)$$

$$\leq \sum_{j\in J} \sum_\ell \hat{q}_\ell x_{\ell,j} \left(1 - e^{-|\ell-j|}\right) + \sum_{j\notin J} \sum_\ell \hat{q}_\ell x_{\ell,j} \left(1 - e^{-|\ell-j|}\right)$$

$$\leq \sum_{j\in J} \sum_\ell \hat{q}_\ell x_{\ell,j} \left(1 - e^{-|\ell-j|}\right) + \sum_{j\notin J} v_j$$

$$\leq \sum_{j\in J} \sum_\ell \hat{q}_\ell x_{\ell,j} \left(1 - e^{-|\ell-j|}\right) + 0.15\sigma$$

36

We get:

$$\sum_{j \in J} \sum_{\ell} \widehat{q}_\ell x_{\ell,j} \left(1 - e^{-|\ell-j|}\right) \geq 0.85\sigma$$

Equivalently, denote $\sum_{j \in J} v_j = \alpha \geq 0.95$.:

$$\sum_{j \in J} (v_j/\alpha) \sum_{\ell} \frac{\widehat{q}_\ell x_{\ell,j}}{v_j/\alpha} \left(1 - e^{-|\ell-j|}\right) \geq 0.85\sigma$$

And by an averaging argument, there exists some $j \in J$ such that $\sum_{\ell} \frac{\widehat{q}_\ell x_{\ell,j}}{v_j/\alpha} \left(1 - e^{-|\ell-j|}\right) \geq 0.85\sigma$. Since $\alpha \geq 0.95$, this is equivalent to:

$$\sum_{\ell} \widehat{q}_\ell x_{\ell,j} \left(1 - e^{-|\ell-j|}\right) \geq 0.8 v_j \sigma$$

$\square$

We next show that for all $j \in J$ we get a strong approximation of the variable $\widetilde{C}_j$ as a function of $s, \widehat{q}_\ell$ and $x_{\ell,j}$:

**Claim 5.5.** *With probability at least* $0.99$ *over the choice of* $(z_1, \ldots, z_s)$, *for every prover response* $(\pi(z_i))_{i \in [s]}$, *characterised by variables* $\{x_{\ell,j}\}_{\ell,j}$ *it holds that:*

$$\mathbb{E}[\widetilde{C}_j] \in \left[ \sum_{\ell:e^{\ell\tau} \geq \tau^2} e^{-\tau} x_{\ell,j} \frac{s^2}{N} \widehat{q}_\ell e^{\ell\tau}, e^{2\tau} \sum_{\ell:e^{\ell\tau} \geq \tau^2} x_{\ell,j} \frac{s^2}{N} \widehat{q}_\ell e^{\ell\tau} \right)$$

*And with probability of at least* $0.99$ *over the choice of* $T$:

$$\left| \widetilde{C}_j - \mathbb{E}_T[\widetilde{C}_j] \right| \leq \tau \mathbb{E}_T[\widetilde{C}_j]$$

The proof of this claim follows Herman and Rothblum [HR22], and is restated in the language of this result in Appendix A.

We thus can conclude that if the run of Tester 5.1.1 ended with accepting the input, then, with high probability over both $(z_i)_{i \in [s]}$ and the samples drawn by $\mathcal{T}$, $\widetilde{C}_j \approx \sum_{\ell} x_{\ell,j} \frac{s^2}{N} \widehat{q}_\ell e^{\ell\tau}$ is close to $\frac{s^2}{N} v_j e^{-j\tau}$. Formally:

**Claim 5.6.** *With probability at least* $0.95$, *if Tester 5.1.1 accepted, then for every* $j \in J$ *it holds that:*

$$\left| \frac{s^2}{N} \sum_{\ell:e^{\ell\tau} \geq \tau^2} x_{\ell,j} \widehat{q}_\ell e^{\ell\tau} - \frac{s^2}{N} v_j e^{j\tau} \right| \leq 10\tau \cdot \frac{s^2}{N} v_j e^{j\tau}$$

*Proof.* From Claim 5.5, we know that with probability at least $0.99$, for every $j \in J$:

$$\left| \widetilde{C}_j - \mathbb{E}\left[\widetilde{C}_j\right] \right| \leq (e^\tau - 1) \mathbb{E}\left[\widetilde{C}_j\right]$$

We conclude that:

$$\left| \widetilde{C}_j - \mathbb{E}\left[\widetilde{C}_j\right] \right| \leq \left(e^{2\tau} - 1\right) \widetilde{C}_j$$

37

And since $\mathbb{E}[\widetilde{C}_j] \in \left[\sum_{\ell:e^{\ell\tau}\geq\tau^2} e^{-\tau}x_{\ell,j}\frac{s^2}{N}\widehat{q}_\ell e^{\ell\tau}, e^{2\tau}\sum_{\ell:e^{\ell\tau}\geq\tau^2} x_{\ell,j}\frac{s^2}{N}\widehat{q}_\ell e^{\ell\tau}\right)$, this implies:

$$\left|\widetilde{C}_j - \frac{s^2}{N}\sum_{\ell:e^{\ell\tau}\geq\tau^2} x_{\ell,j}\widehat{q}_\ell e^{\ell\tau}\right| \leq \left(e^{4\tau}-1\right)\widetilde{C}_j$$

If the tester accepted, it also holds that for the same set of indices $j$ that:

$$\left|\widetilde{C}_j - \frac{s^2}{N}v_j e^{j\tau}\right| \leq 5\tau\frac{s^2}{N}v_j e^{j\tau}$$

And through the triangle inequality, we conclude:

$$\left|\sum_\ell x_{\ell,j}\frac{s^2}{N}\widehat{q}_\ell e^{\ell\tau} - \frac{s^2}{N}v_j e^{j\tau}\right| \leq (e^{4\tau}-1)\widetilde{C}_j + 5\tau\frac{s^2}{N}v_j e^{j\tau} \tag{26}$$

$$\leq \left(e^{4\tau}-1\right)e^\tau\frac{s^2}{N}v_j e^{j\tau} + 5\tau\frac{s^2}{N}v_j e^{j\tau} \tag{27}$$

$$\leq 10\tau\cdot\frac{s^2}{N}v_j e^{j\tau} \tag{28}$$

Where the last inequality stems from the assumptions that $\tau < 0.1$. Dividing by $\frac{s^2}{N}$ we get that if Tester 5.1.1 accepted, then, with probability at least 0.95 over $(z_1,\ldots,z_s)$ and the samples drawn by the tester, for every $j \in J$:

$$\left|\sum_{\ell:e^{\ell\tau}\geq\tau^2}\frac{s^2}{N}\widehat{q}_\ell x_{\ell,j}e^{\ell\tau} - \frac{s^2}{N}v_j e^{j\tau}\right| \leq 10\tau\cdot\frac{s^2}{N}v_j e^{j\tau}$$

$\square$

We now proceed to show that this relation, namely, $\sum_\ell x_{\ell,j}\frac{s^2}{N}\widehat{q}_\ell e^{\ell\tau} \approx \frac{s^2}{N}v_j e^{j\tau}$, which is deduced from the collision matching test assuming the verifier accepted, can be used to show that the alleged size of the $j$'th bucket, expressed as $\frac{v_j}{e^{j\tau}/N}$ is in fact smaller than the expression $\sum_\ell \frac{\widehat{q}_\ell}{e^{\ell\tau}/N}\cdot x_{\ell,j}$.

In order to show this, we first argue the following lemma:

**Lemma 5.7.** *For every* $j \in J$, *if* $\left|v_j e^{j\tau} - \sum_{\ell:e^{\ell\tau}\geq\tau^2}\widehat{q}_\ell x_{\ell,j}e^{\ell\tau}\right| \leq \gamma\cdot v_j e^{j\tau}$, *then, there exists a function* $m_j(\ell)$ *such that:* $m_j(\ell) \in \left[\frac{e^{\ell\tau}}{N}, \frac{e^{j\tau}}{N}\right] \cup \left[\frac{e^{j\tau}}{N}, \frac{e^{\ell\tau}}{N}\right]$, *and:*

$$v_j \log\left(\frac{\sum_{\ell:e^{\ell\tau}\geq\tau^2}\widehat{q}_\ell x_{\ell,j}e^{-\ell\tau}}{v_j e^{-j\tau}}\right) \geq \frac{1}{2}\sum_{\ell:e^{\ell\tau}\geq\tau^2}\widehat{q}_\ell x_{\ell,j}\frac{1}{m_j(\ell)^2}\left(\frac{e^{\ell\tau}}{N}-\frac{e^{j\tau}}{N}\right)^2 - v_j(\gamma+3\tau)$$

*Proof.* Fix $j \in J$. By Taylor's theorem, for every $x \in \mathbb{R}_{\geq 0}$:

$$\log\left(\frac{1}{x}\right) = \log\left(\frac{N}{e^{j\tau}}\right) - \frac{N}{e^{j\tau}}\left(x-\frac{e^{j\tau}}{N}\right) + \frac{1}{2}\cdot\frac{1}{(m(e^{j\tau}/N,x))^2}\left(x-\frac{e^{j\tau}}{N}\right)^2$$

Where $m(e^{j\tau}/N,x) \in \left[x,\frac{e^{j\tau}}{N}\right] \cup \left[\frac{e^{j\tau}}{N},x\right]$. Set $m_j(\ell) = m(e^{j\tau}/N, e^{\ell\tau}/N)$. We get:

$$\sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \log\left(\frac{N}{e^{\ell\tau}}\right) = \sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \log\left(\frac{N}{e^{j\tau}}\right) - \sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \frac{N}{e^{j\tau}} \left(\frac{e^{\ell\tau}}{N} - \frac{e^{j\tau}}{N}\right)$$

$$+ \frac{1}{2} \sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \frac{1}{(m_j(\ell))^2} \left(\frac{e^{\ell\tau}}{N} - \frac{e^{j\tau}}{N}\right)^2$$

Equivalently:

$$\sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \log\left(\frac{N}{e^{\ell\tau}}\right) - \sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \log\left(\frac{N}{e^{j\tau}}\right) = - \sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \frac{N}{e^{j\tau}} \left(\frac{e^{\ell\tau}}{N} - \frac{e^{j\tau}}{N}\right) \tag{29}$$

$$+ \frac{1}{2} \sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \frac{1}{(m_j(\ell))^2} \left(\frac{e^{\ell\tau}}{N} - \frac{e^{j\tau}}{N}\right)^2 \tag{30}$$

First, consider the expression on the left-hand side $\sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \log\left(\frac{N}{e^{\ell\tau}}\right) - \sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \log\left(\frac{N}{e^{j\tau}}\right)$.
Denote $\alpha_j = \sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} = v_j - \widehat{q}_L x_{L,j}$. Recall that by definition, for all $j \in \bar{J}$, $\alpha_j \geq (1-\tau)v_j$.
Therefore, through Jensen's Inequality and the concavity of the log function:

$$\sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \log\left(\frac{N}{e^{\ell\tau}}\right) - \sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \log\left(\frac{N}{e^{j\tau}}\right) = \alpha_j \sum_{\ell:e^{\ell\tau}\geq\tau^2} \frac{\widehat{q}_\ell x_{\ell,j}}{\alpha_j} \log\left(\frac{N}{e^{\ell\tau}}\right) - \alpha_j \log\left(\frac{N}{e^{j\tau}}\right)$$
$$\tag{31}$$

$$\leq \alpha_j \log\left(\frac{\sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} N e^{-\ell\tau}}{\alpha_j}\right) - \alpha_j \log\left(\frac{N}{e^{j\tau}}\right) \tag{32}$$

$$= \alpha_j \log\left(\frac{\sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau}}{\alpha_j e^{-j\tau}}\right) \tag{33}$$

$$\leq v_j \log\left(\frac{\sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau}}{(1-\tau)v_j e^{-j\tau}}\right) \tag{34}$$

$$\leq v_j \log\left(\frac{\sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau}}{v_j e^{-j\tau}}\right) + 2\tau v_j \tag{35}$$

Next, turning our attention to the right-hand side. From the fact that $\sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \geq (1-\tau)v_j$, as well as the assumption $\left|v_j e^{j\tau} - \sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} e^{\ell\tau}\right| \leq \gamma \cdot v_j e^{j\tau}$:

$$\left|\sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \frac{N}{e^{j\tau}} \left(\frac{e^{\ell\tau}}{N} - \frac{e^{j\tau}}{N}\right)\right| = \left|\frac{N}{e^{j\tau}} \left(\sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \frac{e^{\ell\tau}}{N} - \sum_{\ell:e^{\ell\tau}\geq\tau^2} \widehat{q}_\ell x_{\ell,j} \frac{e^{j\tau}}{N}\right)\right| \tag{36}$$

$$= \frac{N}{e^{j\tau}} \left|\left(\sum_\ell \widehat{q}_\ell x_{\ell,j} \frac{e^{\ell\tau}}{N}\right) - v_j(1-\tau)\frac{e^{j\tau}}{N}\right| \tag{37}$$

$$\leq \gamma \cdot v_j + \tau v_j \tag{38}$$

$$\leq (\gamma + \tau)v_j \tag{39}$$

39

Plugging the above inequality, as well as Inequality (31) into Equation (29), we get:

$$v_j \log \left( \frac{\sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau}}{v_j e^{-j\tau}} \right) + 2\tau v_j \geq \frac{1}{2} \sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} \frac{1}{(m_j(\ell))^2} \left( \frac{e^{\ell\tau}}{N} - \frac{e^{j\tau}}{N} \right)^2 - v_j \left( \gamma + \tau \right) \quad (40)$$

Rearranging, we get that for every $j \in J$ for which $\left| v_j e^{j\tau} - \sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{\ell\tau} \right| \leq \gamma \cdot v_j e^{j\tau}$, we get:

$$v_j \log \left( \frac{\sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau}}{v_j e^{-j\tau}} \right) \geq \frac{1}{2} \sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} \frac{1}{(m_j(\ell))^2} \left( \frac{e^{\ell\tau}}{N} - \frac{e^{j\tau}}{N} \right)^2 - v_j \left( \gamma + 3\tau \right) \quad (41)$$

$\square$

Recapping, we proved that with high probability over the choice of samples $S$ and $T$, if Tester 5.1.1 accepted, for every bucket index $j \in J$, it holds that the difference between the alleged size of bucket $j$, $\left( Nv_j e^{-j\tau} \right)$, and the expression $\left( N \sum_\ell \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau} \right)$, depends on $\gamma \cdot v_j$ and $\frac{1}{2} \sum_\ell \widehat{q}_\ell x_{\ell,j} \frac{1}{m_j(\ell)^2} \left( \frac{e^{\ell\tau}}{N} - \frac{e^{j\tau}}{N} \right)^2$; as well as that there exists a $j_0 \in J$ for which $\sum_\ell \widehat{q}_\ell x_{\ell,j_0} \left( 1 - e^{-|\ell - j_0|} \right) > 0.8 v_j \sigma$.

We combine both these ideas to to bound the expression $\sum_\ell \widehat{q}_\ell x_{\ell,j_0} \frac{1}{m_{j_0}(\ell)^2} \left( \frac{e^{\ell\tau}}{N} - \frac{e^{j_0\tau}}{N} \right)^2$:

**Claim 5.8.** *Let* $\sigma \in (0, 0.01)$ *be such that* $\tau \leq \sigma^2$. *Assume* $\frac{1}{s} \sum_{i \in [s]} \left( 1 - \min \left\{ \frac{\pi(z_i)}{D(z_i)}, \frac{D(z_i)}{\pi(z_i)} \right\} \right) > \sigma$, *then there exists some* $j_0 \in J$ *for which:*

$$\sum_\ell \widehat{q}_\ell x_{\ell,j_0} \frac{1}{m_{j_0}(\ell)^2} \left( \frac{e^{\ell\tau}}{N} - \frac{e^{j_0\tau}}{N} \right)^2 \geq 0.5 v_{j_0} \sigma^2$$

*Proof.* By Claim 5.4, with probability at least 0.99 there exists some $j_0 \in J$ for which:

$$\sum_{\ell:e^{\ell\tau} \geq \tau^2} \sum_\ell \widehat{q}_\ell x_{\ell,j_0} \left( 1 - e^{-|\ell - j_0|} \right) > 0.8 v_{j_0} \sigma$$

For every $j \in J$, by definition of $m_j(\ell)$, it holds that:

$$\sum_\ell \widehat{q}_\ell x_{\ell,j} \frac{1}{m_j(\ell)^2} \left( \frac{e^{\ell\tau}}{N} - \frac{e^{j\tau}}{N} \right)^2 \geq \sum_\ell \widehat{q}_\ell x_{\ell,j} \min \left\{ \left( \frac{N}{e^{\ell\tau}} \right)^2, \left( \frac{N}{e^{j\tau}} \right)^2 \right\} \left( \frac{e^{\ell\tau}}{N} - \frac{e^{j\tau}}{N} \right)^2 \quad (42)$$

$$\geq \sum_\ell \widehat{q}_\ell x_{\ell,j} \left( 1 - e^{-|\ell - j|} \right)^2 \quad (43)$$

In particular, for $j_0$, since $j_0 \in J$, it holds that $\sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j_0} = \alpha_{j_0} \geq (1 - \tau) v_{j_0}$, and through Jensen's Inequality:

$$\sum_\ell \widehat{q}_\ell x_{\ell,j_0} \left( 1 - e^{-|\ell - j_0|} \right)^2 = \alpha_{j_0} \sum_\ell \frac{\widehat{q}_\ell x_{\ell,j_0}}{\alpha_{j_0}} \left( 1 - e^{-|\ell - j_0|} \right)^2$$

$$\geq \alpha_{j_0} \left( \sum_\ell \frac{\widehat{q}_\ell x_{\ell,j_0}}{\alpha_{j_0}} \left( 1 - e^{-|\ell - j_0|} \right) \right)^2$$

$$\geq 0.6 \alpha_{j_0} \sigma^2$$

$$\geq 0.5 v_{j_0} \sigma^2$$

$\square$

40

**Proposition 5.9** (Soundness of Tester 5.1.1)**.** *With high probability over the choice of* $(z_i)_{i \in [s]}$, *for every set of claims* $(\pi(z_i))_{i \in [s]}$, *if* $\mathcal{T}$ *accepted, then, with probability of at least* 0.95 *over the samples drawn by* $\mathcal{T}$, *if* $\sum_{i \in [s]} \left( 1 - \min \left\{ \frac{\pi(z_i)}{D(z_i)}, \frac{D(z_i)}{\pi(z_i)} \right\} \right) \geq \sigma$ *it holds that there exists some* $j$ *such that* $v_j \geq \frac{\tau^2}{\log N}$, $\frac{e^{j\tau}}{N} \geq \frac{\tau}{N}$ *and:*

$$\sum_{\ell : e^{\ell \tau} \geq \tau^2} \widehat{q}_\ell x_{\ell, j} e^{-\ell \tau} \geq \left( 1 + 0.2\sigma^2 - 10\tau \right) v_j e^{-j\tau}$$

*Proof.* By Claim 5.6, if the tester accepted, then with probability at least 0.9, it holds that for every $j \in J$ $\left| \frac{s^2}{N} \sum_\ell x_{\ell,j} \widehat{q}_\ell e^{\ell \tau} - \frac{s^2}{N} v_j e^{j\tau} \right| \leq 10\tau \cdot \frac{s^2}{N} v_j e^{j\tau}$.

Through Lemma 5.7, we conclude that for all such $j$, plugging $\gamma = 10\tau$, it also holds that there exists some $m_j(\ell) \in \left[ \frac{e^{\ell \tau}}{N}, \frac{e^{j\tau}}{N} \right] \cup \left[ \frac{e^{j\tau}}{N}, \frac{e^{\ell \tau}}{N} \right]$, for which:

$$v_j \log \left( \frac{\sum_{\ell : e^{\ell \tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell \tau}}{v_j e^{-j\tau}} \right) \geq \frac{1}{2} \sum_{\ell : e^{\ell \tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} \frac{1}{m_j(\ell)^2} \left( \frac{e^{\ell \tau}}{N} - \frac{e^{j\tau}}{N} \right)^2 - 10\tau v_j$$

Finally, by Claim 5.8, there must be some $j \in J$ for which:

$$v_j \log \left( \frac{\sum_{\ell : e^{\ell \tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell \tau}}{v_j e^{-j\tau}} \right) \geq 0.2 v_j \sigma^2 - 10\tau v_j$$

Rearranging the last inequality, we get:

$$\sum_{\ell : e^{\ell \tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell \tau} \geq e^{0.2\sigma^2 - 10\tau} v_j e^{-j\tau} \geq \left( 1 + 0.2\sigma^2 - 10\tau \right) v_j e^{-j\tau}$$

$\square$

# 6  Full Protocol

In this section we prove the Theorem 1.1: a doubly-efficient proof system for any label-invariant property. The full construction ties together the *Collision-Matching Tester* of Section 5 and the *Bucket Size Verification Protocol* of Section 4. Tying these two components together requires that the distribution satisfy the following promise: there should be no "heavy" elements (whose probabilities are above $\frac{\tau}{\sqrt{N}}$ for the chosen accuracy parameter $\tau$, see above). Note that a distribution with only *heavy* elements can be easily learned (i.e. obtaining an explicit description of a distribution $D'$ that is arbitrarily close to $D$) by the verifier without any interaction and low sample complexity (roughly $O(\sqrt{N})$ samples, see Theorem 3.15), while a distribution solely supported on *light* elements (i.e. with probability smaller than $\frac{\tau}{\sqrt{N}}$) requires significantly more samples (potentially $\Omega(N)$) to be learned without interaction. Therefore, the main focal point of our construction deals with verifying properties with distribution that are *harder* to learn, and assumes that no heavy elements are found in the support. Nonetheless, a general distribution might have both light and heavy elements. In this section we address this difficulty, as well as put together all the tools presented in previous sections.

The full construction proceeds in two steps:

1. In Section 6.1, we present a doubly-efficient protocol for reconstructing the histogram of a distribution without heavy elements, combining the tools presented in the previous sections, putting together the protocol of Section 4 and the collision-matching test of Section 5.

2. In Section 6.2, we show how to obtain an approximate histogram of a general distribution. We recall a reduction of [HR22] that allows us to essentially "get rid" of any heavy elements: in a nutshell, the verifier can use distribution learning to identify *all the heavy elements* and approximate their probabilities. Then, if these elements account for less than $1-\sigma$ of the mass (for some distance parameter $\sigma \in (0, 0.1)$), then there is enough mass on *light elements*. The verifier then runs the protocol of Section 6 on the light-elements through rejection sampling, and combines both histograms together.

3. In Section 6.3 we refer again to [HR22] who show how to leverage a verified approximate histogram of $D$ to a proof system for label-invariant distribution properties given sample access to $D$.

## 6.1  Handling Light Elements

**Theorem 6.1.** *There exists a 2-message interactive protocol between an honest verifier and a (potentially malicious) prover, where the verifier receives as input parameters $\sigma \in (0, 0.1)$ and $100 < N \in \mathbb{N}$, as well as sample access to a distribution $D$ over domain $[N]$. Set $\tau = \frac{1}{500}\sigma^2$. Assume $D(x) \leq \frac{\tau}{\sqrt{N}}$. The communication complexity, verifier sample complexity, and verifier runtime are all $s = \widetilde{O}(\sqrt{N}) \cdot \mathsf{poly}(\sigma^{-1})$. Given sample access to the distribution $D$, the honest prover requires with high probability $\widetilde{O}(N)\,\mathsf{poly}(\sigma^{-1})$ samples and runtime.*

*At the end of the interaction, the verifier rejects or outputs $((z_i, \pi(z_i)))_{i \in [s]}$, where $(z_i)_{i \in [s]}$ is a sample of size $s$ drawn i.i.d. by $D$, such that:*

- *If the prover is honest, then with probability at least $0.9$, the verifier doesn't reject, and $((z_i, \pi(z_i)))_{i \in [s]}$ satsifies $\frac{1}{s}\sum_{i \in [s]}\left(1 - \min\left\{\frac{\pi(z_i)}{D(z_i)}, \frac{D(z_i)}{\pi(z_i)}\right\}\right) = O(\tau)$.*

- *Whatever strategy a dishonest prover follows, with probability at most* $0.1$ *over the verifier's coin tosses and samples, they accept and* $((z_i, \pi(z_i)))_{i \in [s]}$ *satisfies:*

$$\frac{1}{s} \sum_{i \in [s]} \left(1 - \min\left\{\frac{\pi(z_i)}{D(z_i)}, \frac{D(z_i)}{\pi(z_i)}\right\}\right) \geq \sigma$$

---

**Protocol 6.1.1: Doubly-Efficient Histogram Retrieval Protocol**

**Input:** parameters $N \in \mathbb{N}$, $\sigma \in (0,1)$, as well as sample access to distribution $D$ over domain $[N]$. $D$ is assumed to satisfy $\forall x$, $D(x) \leq \frac{\tau}{\sqrt{N}}$ for $\tau = \frac{1}{500}\sigma^2$.

**Goal:** obtain with high probability $((z_i, \pi(z_i)))_{i \in [s]}$, such that $(z_i)_{i \in [s]}$ is a sample drawn i.i.d. by $D$, and $\frac{1}{s} \sum_{i \in [s]} \left(1 - \min\left\{\frac{\pi(z_i)}{D(z_i)}, \frac{D(z_i)}{\pi(z_i)}\right\}\right) \leq \sigma$, for $s = \widetilde{O}(\sqrt{N})\mathsf{poly}\left(\sigma^{-1}\right)$.

1. V-P: run Protocol 4.1.1 with the distribution $D$, parameter $\tau$, and $s = \widetilde{\theta}\left(\sqrt{N}\right) \cdot \mathsf{poly}(\sigma^{-1})$. If protocol ended without rejection, obtain a tagged sample $((z_i, \pi(z_i)))_{i \in [s]}$ for some $s = \widetilde{\theta}\left(\sqrt{N}\right) \cdot \mathsf{poly}(\sigma^{-1})$.

2. V: Run Tester 5.1.1 on the obtained tagged sample $((z_i, \pi(z_i)))_{i \in [s]}$, with parameters $N, \sigma, \tau$. Reject if tester rejected.

3. V: output $((z_i, \pi(z_i)))_{i \in [s]}$

---

### 6.1.1 Completeness of Protocol 6.1.1

Assume the prover is honest. We can conclude the following:

**Step (1) of Protocol 6.1.1.** After running Step (1) of Protocol 6.1.1, the verifier obtains a tagged sample $((z_i, \pi(z_i)))_{i \in [s]}$, from which they can deduce a $(N, \tau)$-histogram $\{v_j\}_j$, as given in Protocol 4.1.1. If the prover is honest, then by the completeness condition of Protocol 4.1.1 outlined in Proposition 4.1.1, at the end of the step it holds that with probability at least $0.95$:

- The verifier accepts.

- For every $i \in [s]$, $\pi(z_i) \in D(z_i)\left[e^{-\tau}, e^{\tau}\right]$, and so: $\frac{1}{s} \sum_{i \in [s]} \left(1 - \min\left\{\frac{\pi(z_i)}{D(z_i)}, \frac{D(z_i)}{\pi(z_i)}\right\}\right) = O(\tau) = O(\sigma^2)$

**Step (2) of Protocol 6.1.1** . Assuming the conditions above hold, then the input to Tester 5.1.1 is a correctly tag sample $S = ((z_i, \mathsf{tag}(z_i)))_{i \in [s]}$, with $s = \widetilde{O}\left(\sqrt{N}\right)\mathsf{poly}(\sigma^{-1})$, and by the completeness condition of Tester 5.1.1, outlined in Proposition 5.1, with probability at least $0.99$, the Tester accepts.

### 6.1.2 Soundness of Protocol 6.1.1

In order to show soundness, we assume that Step (1) passed without rejection, and that the verifier obtained a tagged sample $((z_i, \pi(z_i)))_{i \in [s]}$ such that $\frac{1}{s} \sum_{i \in [s]} \left( 1 - \min \left\{ \frac{\pi(z_i)}{D(z_i)}, \frac{D(z_i)}{\pi(z_i)} \right\} \right) \geq \sigma$. We proceed to show that with high probability Step (2) results in rejection.

To understand how we achieve this, we first review what we know after having passed Step (1). With probability at least $0.95$ over the randomness of $V$, based on the soundness of Protocol 4.1.1 outlined in Section 4 we are guaranteed the following: define $\{v_j\}_j$ as well as $\{x_{\ell,j}\}_{\ell,j}$ as in Section 4. Every bucket index $j$ such that $\frac{e^{j\tau}}{N} \geq \frac{\tau}{N}$ and $v_j \geq \frac{\tau^2}{\log N}$ satisfies:

$$v_j e^{-j\tau} \geq (1 - 18\tau) \sum_{\ell : e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell e^{-\ell\tau} x_{\ell,j} - f\left( \log N, \tau^{-1} \right) \cdot \frac{1}{s} \tag{44}$$

Where $f\left( \log N, \tau^{-1} \right) = (\log N)^{t_1} \cdot \left( \frac{1}{\tau} \right)^{t_2}$, for some $t_1, t_2 \in \mathbb{N}$. If we assume that the tester at Step (2) passed then we are also guaranteed that with high probability there exists some bucket index $j$ such that $e^{j\tau} \geq \tau$, $v_j \geq \frac{\tau^2}{\log N}$ and:

$$\sum_{\ell : e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau} \geq \left( 1 + 0.2\sigma^2 - 10\tau \right) v_j e^{-j\tau} \tag{45}$$

That means that there exists some bucket $j$ for which:

$$(1 - 18\tau) \sum_{\ell : e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell e^{-\ell\tau} x_{\ell,j} - f\left( \log N, \tau^{-1} \right) \cdot \frac{1}{s} \leq v_j e^{-j\tau} \leq \frac{1}{(1 + 0.2\sigma^2 - 10\tau)} \sum_{\ell : e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau} \tag{46}$$

We show that assuming $\frac{1}{s} \sum_{i \in [s]} \left( 1 - \min \left\{ \frac{\pi(z_i)}{D(z_i)}, \frac{D(z_i)}{\pi(z_i)} \right\} \right) \geq \sigma$, we can choose $s = \widetilde{\theta}\left( \sqrt{N} \right) \mathsf{poly}\left( \sigma^{-1} \right)$ for which this cannot be, and so, it cannot be that with high probability both the *Bucket Size Verification Protocol* passes and so does the *Collisions Matching Tester*.

**Claim 6.2.** *There exists a choice of $s = \widetilde{O}\left( \sqrt{N} \right) \mathsf{poly}(\sigma^{-1})$ such that for every $j \in J$:*

$$(1 - 18\tau) \sum_{\ell : e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell e^{-\ell\tau} x_{\ell,j} - f\left( \log N, \tau^{-1} \right) \cdot \frac{1}{s} > \frac{1}{(1 + 0.2\sigma^2 - 10\tau)} \sum_{\ell : e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau} \tag{47}$$

*Proof.* First, note that:

$$(1 - 18\tau) \sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau} - \frac{1}{1 + 0.2\sigma^2 - 10\tau} \sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau}$$

$$= \left(1 - 18\tau - \frac{1}{1 + 0.2\sigma^2 - 10\tau}\right) \sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau}$$

$$\geq \left(1 - 18\tau - \left(1 - 0.1\sigma^2 + 5\tau\right)\right) \sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau}$$

$$\geq \left(0.1\sigma^2 - 23\tau\right) \sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau}$$

$$\geq \left(0.1\sigma^2 - 23\tau\right) \cdot \frac{\sqrt{N}\tau^2}{2 \log N}$$

Where the last inequality stems from:

$$\sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau} \geq \frac{\sqrt{N}}{\tau} \sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} \geq (1 - \tau)v_j \frac{\sqrt{N}}{\tau} \geq \frac{\sqrt{N}\tau^2}{2 \log N}$$

Which is justified by the assumption that $j \in J$, as well as the assumption $q_\ell = 0$ for $\ell$ satisfying $\frac{e^{\ell\tau}}{N} \geq \frac{\tau}{\sqrt{N}}$.

Next, since $\tau = \frac{1}{500}\sigma^2$, we get that:

$$(1 - 18\tau) \sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau} - \frac{1}{1 + 0.2\sigma^2 - 10\tau} \sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{-\ell\tau} \geq \left(0.1\sigma^2 - 23\tau\right) \cdot \frac{\sqrt{N}\tau^2}{2 \log N} \geq 0.05 \frac{\sigma^2\tau^2\sqrt{N}}{2 \log N}$$

Thus, if we set $s$ such that $\frac{1}{s}f\left(\log N, \tau^{-1}\right) \leq 0.05\frac{\sigma^2\tau^2\sqrt{N}}{2 \log N}$, we're done. Indeed, we conclude the proof by setting:

$$s = 50\sqrt{N} \cdot f(\log N, \tau^{-1}) \cdot \frac{\log N}{\sigma^2\tau^2} = \widetilde{O}(\sqrt{N})\mathsf{poly}(\tau^{-1})$$

$\square$

### 6.1.3 On $\Delta_{\mathrm{RL}}$-distance guaranatees

In [HR22] the authors consider the $\Delta_{\mathrm{RL}}$ distance measure between distribution (see Definition 3.9). In short, for two distributions $\Delta_{\mathrm{RL}}(P,Q)$ (called the *relabeling distance between P and Q*) can be thought of as the smallest distance between $P$ (equivalently $Q$) and a permutation of $Q$ (equivalently, $P$), where a permutation of distribution is defined in Definition 3.8.

As we are interested in the distance between $D$ and a label-invariant distribution property, which is a set closed under permutations. Two distributions that are permutation of one another might be very far in $\Delta_{\mathrm{SD}}$, however, if one is close to a label-invariant property, so is the other. And so, when talking about label invariant properties, it is useful to consider *relabeling distance*, $\Delta_{\mathrm{RL}}$.

As a consequence, [HR22] use $\Delta_{\mathrm{RL}}$ distance to characterize the soundness condition of their *histogram reconstruction protocol*. They require that their proof system will reject histograms

$\{a_j\}_j$ for which $\Delta_{\mathrm{RL}}(\{a_j\}_j, D) \geq \sigma$ with high probability. In this work we use a different sense of soundness. We require the verifier to reject with high probability interactions that produce a tagged sample $(z_i, \pi(z_i)))_{i \in [s]}$ that satisfies $\frac{1}{s} \sum_{i \in [s]} \left(1 - \min\left\{\frac{D(z_i)}{\pi(z_i)}, \frac{\pi(z_i)}{D(z_i)}\right\}\right) > \sigma$. This measure is more delicate, as it sums the individual mislabeling *per each sample*, instead of considering the implicit aggregate distance measure of $\Delta_{\mathrm{RL}}$. In this section we show that the new soundness condition is compatible with the one in [HR22]. That is, if $\{a_j\}_j$ was obtained through our protocol, then if $\Delta_{\mathrm{RL}}(D, \{a_j\}_j) \geq \sigma$, then it must be that $\frac{1}{s} \sum_{i \in [s]} \left(1 - \min\left\{\frac{D(z_i)}{\pi(z_i)}, \frac{\pi(z_i)}{D(z_i)}\right\}\right) \geq \sigma/4$. That is, running our protocol with $\sigma/4$ guarantees that every histogram rejected by [HR22] with high probability is also rejected by our protocol, despite having a different soundness condition. This allows us to plug our result in the setting offered by [HR22].

**Remark 6.3.** *The main benefit of considering the soundness condition with respect to the quantity* $\frac{1}{s} \sum_{i \in [s]} \left(1 - \min\left\{\frac{D(z_i)}{\pi(z_i)}, \frac{\pi(z_i)}{D(z_i)}\right\}\right)$ *is that it allows us to argue more delicately about the mistagging, as the prover is penalized per mistagged sample. This is a basis for extending the protocol to non-label invariant properties, where we want to understand how far each sample is from it's true tag.*

**Claim 6.4.** *Fix $((z_i, \pi(z_i)))_{i \in [s]}$ as defined in the previous sections, and let $\{v_j\}_j$ be the $(N, \tau)$-histogram induced by it. If $\Delta_{RL}(\{v_j\}_j, D) \geq \sigma$, then:*

$$\frac{1}{s} \sum_{i \in [s]} \left(1 - \min\left\{\frac{\pi(z_i)}{D(z_i)}, \frac{D(z_i)}{\pi(z_i)}\right\}\right) \geq \sigma/4$$

*Proof.* Through the proof of Claim 5.3, we learn that if:

$$\sum_j \sum_\ell \widehat{q}_\ell x_{\ell,j} \left(1 - e^{-|\ell - j|\tau}\right) \geq \sigma/4 \tag{48}$$

Then:

$$\frac{1}{s} \sum_{i \in [s]} \left(1 - \min\left\{\frac{\pi(z_i)}{D(z_i)}, \frac{D(z_i)}{\pi(z_i)}\right\}\right) \geq \sigma/4$$

Therefore, assuming $\Delta_{\mathrm{RL}}(\{v_j\}_j, D) \geq \sigma$, we show that Inequality (48) holds. First, we show this under the assumption that $\{v_j\}_j$ is a histogram of a (roughly)-uniform distribution. I.e. there exists some $j_0$ for which $v_{j_0} = 1$. This also implies that $x_{\ell,j} = 1$ if $j = j_0$, and $x_{\ell,j} = 0$ otherwise. Then we decompose a general histogram to many instances of this case.

**Step 1. There exists some bucket index $j_0$ s.t. $v_{j_0} = 1$.** In this case, we want to show:

$$\sum_\ell \widehat{q}_\ell x_{\ell,j_0} \left(1 - e^{-|\ell - j_0|\tau}\right) = \sum_\ell \widehat{q}_\ell \left(1 - e^{-|\ell - j_0|\tau}\right) \geq \sigma/5 \tag{49}$$

Through [HR22] it holds that with high probability over $(z_i)_{i \in [s]}$, $\Delta_{\mathrm{RL}}(D, \{\widehat{q}_\ell\}_\ell) \leq 5\tau$, as $\{\widehat{q}_\ell\}_\ell$ represents the empirical histogram of $D$ for a large sample.

Let $D'$ be a distribution over $[N]$ consistent with $(N, \tau)$-histogram $\{\widehat{q}_\ell\}_\ell$, and assume w.l.o.g. that for all $x \in [N - 1]$, $D'(x) \geq D'(x + 1)$. Let $P$ be a distribution consistent with $\{v_j\}_j$ (i.e. $P$ has only one bucket - the $j_0$'th bucket, and is thus almost uniform), such that $\Delta_{\mathrm{RL}}(D', \{v_j\}_j) = \Delta_{\mathrm{SD}}(D', P)$, which also implies that for all $x \in [N]$, $P(x) \geq P(x + 1)$ (as showed in [HR22]).

Define:

$$A_1 = \left\{ i \in \text{Supp}(D') \cap \text{Supp}(P) \ : \ D'(i) \geq P(i) \neq 0 \right\}$$
$$A_2 = \left\{ i \in \text{Supp}(D') \cap \text{Supp}(P) \ : \ D'(i) < P(i) \neq 0 \right\}$$
$$A_3 = \text{Supp}(D') \setminus \text{Supp}(P)$$
$$A_4 = \text{Supp}(P) \setminus \text{Supp}(D')$$

Observe that $\sigma - 5\tau \leq \Delta_{\text{SD}}(P, D') = \sum_{x \in A_1} (D'(x) - P(x)) + \sum_{x \in A_3} D'(x)$. Denote $\sigma' = \sigma - 5\tau$. Thus, it must be that either $\sum_{x \in A_1} (D'(x) - P(x)) \geq \sigma'/2$, or $\sum_{x \in A_3} D'(x) \geq \sigma'/2$. In order to prove that Inequality (49), we consider the following case analysis:

**Case 1.1:** assume $\sum_{x \in A_1} (D'(x) - P(x)) \geq \sigma'/2$ . define:

$$A_{good} = \left\{ i \in A_1 : \frac{D'(i) - \frac{e^{j_0 \tau}}{N}}{D'(i)} > \sigma'/4 \right\}$$

$$A_{bad} = \left\{ i \in A_1 : \frac{D'(i) - \frac{e^{j_0 \tau}}{N}}{D'(i)} \leq \sigma'/4 \right\}$$

Note that:

$$\sum_{i \in A_{bad}} \left( D'(i) - \frac{e^{j_0 \tau}}{N} \right) = \sum_{i \in A_{bad}} D'(i) \frac{\left( D'(i) - \frac{e^{j_0 \tau}}{N} \right)}{D'(i)} \leq \frac{\sigma'}{4} \sum_{i \in A_{bad}} D'(i) \leq \frac{\sigma'}{4}$$

Since $\sum_{i \in A_1} \left( D'(i) - \frac{e^{j_0 \tau}}{N} \right) \geq \sigma'/2$, this implies that:

$$\sum_{i \in A_{good}} \left( D'(i) - e^{j_0 \tau}/N \right) \geq \sigma'/4$$

And so:

$$\sum_\ell \widehat{q}_\ell x_{\ell,j_0} \left(1 - e^{-|\ell-j_0|\tau}\right) \geq \sum_{\ell:\ell\geq j_0} \widehat{q}_\ell \left(1 - e^{-|\ell-j_0|\tau}\right)$$

$$\geq \sum_{\ell:\ell\geq j_0} \widehat{q}_\ell \left(1 - e^{j_0-\ell\tau}\right)$$

$$\geq \sum_{\ell:\ell\geq j_0} \frac{\widehat{q}_\ell}{e^{\ell\tau}/N} \left(\frac{e^{\ell\tau}}{N} - \frac{e^{j_0}}{N}\right)$$

$$\geq (1-2\tau) \sum_{\ell:\ell\geq j_0} \sum_{i\in B_\ell^{D'}} \left(\frac{e^{\ell\tau}}{N} - \frac{e^{j_0}}{N}\right)$$

$$\geq (1-2\tau)^2 \sum_{i\in A_1} \left(D'(i) - \frac{e^{j_0}}{N}\right)$$

$$\geq (1-2\tau)^2 \sum_{i\in A_{good}} \left(D'(i) - \frac{e^{j_0}}{N}\right)$$

$$\geq (1-2\tau)^2 \sigma'/4$$

$$\geq \sigma/5$$

**Case 1.2: assume** $\sum_{x\in A_3} D'(i) \geq \sigma'/2$. In particular it holds that $A_4 = \phi$. This second case is divided into two subcases, according to the value of $p_{min}^{A_2} = \min\{D'(i) : i \in A_2\}$.

**Case 1.2.1: assume** $p_{min}^{A_2} < \frac{e^{j_0\tau}}{2N}$, **and** $\sum_{i\in A_3} D'(i) \geq \sigma'/2$. Recall that we assumed without loss of generality that $D'(i) \geq D'(i+1)$ as well as $P(i) \geq P(i+1)$, for all $i \in [N]$. Therefore, in particular, we conclude that for every $i \in A_3$, $D'(i) \leq p_{min}^{A_2}$ - this is justified by observing that $A_2 \subseteq \text{Supp}(P)$, while $A_3 \cap \text{Supp}(P) = \phi$, and so, we deduce that for all $j \in A_3$ and $k \in A_2$, $j > k$. Therefore, for every $i \in A_3$, $\left(e^{j_0\tau}/N - D'(i)\right) \geq \frac{e^{j_0\tau}}{2N}$, which implies: $\frac{(e^{j_0\tau}/N - D'(i))}{D'(i)} \geq \frac{e^{j_0\tau}/2N}{D'(i)} \geq 1$.

And so:

$$\sum_\ell \widehat{q}_\ell x_{\ell,j_0} \left(1 - e^{-|\ell - j_0|\tau}\right) \geq \sum_{\ell : \frac{e^{\ell\tau}}{N} \leq \frac{e^{j_0\tau}}{2N}} \widehat{q}_\ell \left(1 - e^{(\ell - j_0)\tau}\right)$$

$$\geq \sum_{\ell : \frac{e^{\ell\tau}}{N} \leq \frac{e^{j_0\tau}}{2N}} \frac{\widehat{q}_\ell}{e^{j_0\tau}/N} \left(\frac{e^{j_0}}{N} - \frac{e^{\ell\tau}}{N}\right)$$

$$= \sum_{\ell : \frac{e^{\ell\tau}}{N} \leq \frac{e^{j_0\tau}}{2N}} \frac{e^{\ell\tau}/N}{e^{j_0\tau}/N} \cdot \frac{\widehat{q}_\ell}{e^{\ell\tau}/N} \left(\frac{e^{j_0}}{N} - \frac{e^{\ell\tau}}{N}\right)$$

$$\geq (1 - 2\tau) \sum_{\ell : \frac{e^{\ell\tau}}{N} \leq \frac{e^{j_0\tau}}{2N}} \frac{e^{\ell\tau}/N}{e^{j_0\tau}/N} \cdot \left|B_\ell^{D'}\right| \left(\frac{e^{j_0}}{N} - \frac{e^{\ell\tau}}{N}\right)$$

$$\geq (1 - 2\tau)^2 \sum_{i \in A_3} \frac{D'(i)}{e^{j_0\tau}/N} \cdot \left(\frac{e^{j_0}}{N} - D'(i)\right)$$

$$\geq (1 - 2\tau)^2 \sum_{i \in A_3} \frac{D'(i)}{e^{j_0\tau}/N} \cdot \frac{e^{j_0}}{2N}$$

$$\geq (1 - 2\tau)^2 \sum_{i \in A_3} \frac{D'(i)}{2}$$

$$\geq \sigma/5$$

**Case 1.2.2: Assume** $p_{min}^{A_2} \geq \frac{e^{j_0\tau}}{2N}$ **and** $\sum_{i \in A_3} D'(i) \geq \sigma'/2$ . This implies that $A_4 = \phi$. We have that $\Delta_{\mathrm{SD}}(D', P) = \sum_{i \in A_2} (P(i) - D'(i))$:

$$\sum_\ell \widehat{q}_\ell x_{\ell,j_0} \left(1 - e^{-|\ell - j_0|\tau}\right) \geq \sum_{\ell : \frac{e^{\ell\tau}}{N} \in \left[\frac{e^{j_0\tau}}{2N}, \frac{e^{j_0\tau}}{N}\right]} \widehat{q}_\ell \left(1 - e^{(\ell - j_0)\tau}\right)$$

$$\geq \sum_{\ell : \frac{e^{\ell\tau}}{N} \in \left[\frac{e^{j_0\tau}}{2N}, \frac{e^{j_0\tau}}{N}\right]} \frac{\widehat{q}_\ell}{e^{j_0\tau}/N} \left(\frac{e^{j_0}}{N} - \frac{e^{\ell\tau}}{N}\right)$$

$$\geq \sum_{\ell : \frac{e^{\ell\tau}}{N} \in \left[\frac{e^{j_0\tau}}{2N}, \frac{e^{j_0\tau}}{N}\right]} \cdot \frac{1}{2} \frac{\widehat{q}_\ell}{e^{\ell\tau}/N} \left(\frac{e^{j_0}}{N} - \frac{e^{\ell\tau}}{N}\right)$$

$$\geq \frac{(1 - 2\tau)}{2} \sum_{\ell : \frac{e^{\ell\tau}}{N} \in \left[\frac{e^{j_0\tau}}{2N}, \frac{e^{j_0\tau}}{N}\right]} \left|B_\ell^{D'}\right| \left(\frac{e^{j_0}}{N} - \frac{e^{\ell\tau}}{N}\right)$$

$$\geq \frac{(1 - 2\tau)^2}{2} \sum_{i \in A_2} \left(\frac{e^{j_0}}{N} - D'(i)\right)$$

$$\geq \sigma/5$$

And we conclude that in any case, $\sum_\ell \widehat{q}_\ell x_{\ell,j_0} \left(1 - e^{-|\ell - j_0|\tau}\right) \geq \sigma/5$. We now prove the general case:

**Case 2: $\{v_j\}_j$ is supported on multiple buckets.** We rewrite Inequality (48) we wish to prove as follows:

$$\sum_j v_j \sum_\ell \frac{\widehat{q}_\ell x_{\ell,j}}{v_j} \left(1 - e^{-|\ell - j|\tau}\right) \geq \sigma/5$$

Consider thus a new distribution $P$ consistent with $\{v_j\}_j$ that's defined as follows: there exists a subdomain $X_j \subseteq [N]$ such that $D'|_{X_j}$ has the conditional histogram $\left\{\frac{\widehat{q}_\ell x_{\ell,j}}{v_j}\right\}_{f(\ell,j)}$ Where for every $j, \ell$, define $f(\ell, j) = \log\left(N \cdot \frac{e^{\ell\tau}}{Nv_j}\right)/\tau$, and $P|_{X_j}$ is uniform over $\frac{v_j}{e^{j\tau}/N}$ elements, and also satisfies $\Delta_{\mathrm{RL}}(D'|_{X_j}, P|_{X_j}) = \Delta_{\mathrm{SD}}(D'|_{X_j}, P|_{X_j}) = \sigma_j$. From the above analysis, we get that for every $\ell$:

$$\sum_\ell \frac{\widehat{q}_\ell x_{\ell,j}}{v_j} \left(1 - e^{-|\ell - j|\tau}\right) \geq \sigma_j/5$$

And so, taking weighted sum according $v_j$ on both sides of the inequality, we get:

$$\sum_j v_j \sum_\ell \frac{\widehat{q}_\ell x_{\ell,j}}{v_j} \left(1 - e^{-|\ell - j|\tau}\right) \geq \sum_j v_j\sigma_j/5$$

And since $\sum_j v_j \Delta_{\mathrm{SD}}\left(D|_{X_j}, P|_{X_j}\right) \geq \sum_j v_j\sigma_j/4 \geq \sigma/4$, we get:

$$\sum_j v_j \sum_\ell \frac{\widehat{q}_\ell x_{\ell,j}}{v_j} \left(1 - e^{-|\ell - j|\tau}\right) \geq \sum_j v_j\sigma_j/5 = \sigma/5$$

$\square$

## 6.2 Handling General Distributions

In the previous section we showed how to handle distributions with no *heavy* elements. In this section we demonstrate how to decompose a general distribution's support into *light* and *heavy* subdomains, how to approximate the conditional histograms of $D$ for both sections, and then how to combine them for a complete histogram of $D$. We do so by following [HR22], and so, we only bring here this procedure in broad strokes. We refer the reader to the *IP for verified histogram reconstruction* from [HR22] for further details.

### 6.2.1 Short overview of [HR22] IP for verified histogram reconstruction

Given distance parameter $\sigma$, and accuracy parameter $\tau \leq \frac{1}{500}\sigma^2$, the verifier divides the domain into two subsets $X_{light}$ and $X_{heavy}$, obtains an approximate $(N, \tau)$-histogram of $D$ restricted to these subdomains, and then proceeds to combine the results obtained on each subdomain into one verified histogram for $D$, $\{a_j\}_j$, such that if the prover is honest, with high probability the verifier accepts and $\Delta_{\mathrm{RL}}(D, \{a_j\}_j) = O(\tau)$; and no matter how the prover tries to cheat, at the end of the interaction, with high probability, either the verifier rejects or $\Delta_{\mathrm{RL}}(D, \{a_j\}_j) = O(\sigma^2)$.

**Indentifying the heavy elements.** The verifier, on its own (i.e. without interacting with the prover), draws sufficiently many samples to identify a set $W_D = \left\{ x \in [N] : D(x) \geq \frac{\tau}{\sqrt{N}} \right\}$ of bounded size.

This is simply done by considering the set $X_{heavy}$ to be the set of distinct elements observed after taking $\widetilde{O}(\sqrt{N}/\tau^2)$ samples from $D$. W.h.p. all elements with probability at least $\frac{\tau}{\sqrt{N}}$ will appear in this set at least once, i.e. $W_D \subseteq X_{heavy}$ (note that $X_{heavy}$ might also contain light elements, this is fine, the focus is to capture all the *heavy* elements). Let $X_{light}$ be the complement, i.e. $([N] \setminus X_{heavy})$. We can estimate the probabilities of the sets $X_{heavy}$ and $X_{light}$ by $D$ up to $\sigma/100$ accuracy using $\mathsf{poly}(1/\sigma)$ samples. Note that if distribution $D$ has no or very few heavy elements, we'll get that the estimate $D(X_{heavy})$ is very small.

**Finding the histogram of $D$ conditioned on both the *heavy* subdomain and the *light* subdomain.** We now proceed with a case analysis: we'd like to obtain verified histograms of both $D\big|_{X_{heavy}}$ and $D\big|_{X_{light}}$, after having identified $X_{heavy}$ as explained above. If both sets have weight (by $D$) that is sufficiently bounded away from 0, then we can do so by obtain sufficiently many samples from each of them (and, for $X_{light}$, interacting with the prover). If, on the other hand, one of the two sets has tiny probability, then we might not be able to obtain sufficiently many samples from it, but on the other hand it will not affect $D$'s histogram (because the set's total probability is small), and we can safely ignore it. We need only verify the histogram of the other set (from which we can sample, since its probability is close to 1).

For $X_{heavy}$, if our estimate on its probability is less than least $\sigma/9$, then we can safely ignore it. Otherwise, we learn its histogram, and observe that in this case, we know that $X_{heavy}$'s true probability by $D$ is at least $\sigma/10$, and so we can learn its histogram up to a $\sigma/10$ statistical distance using the *folklore distribution learner* (see Section 3).[4] We can sample from $D\big|_{X_{heavy}}$ using rejection sampling, since its weight is sufficiently large. Thus, using also the fact that $|X_{heavy}| = \widetilde{O}\left(\sqrt{N}/\tau\right)$, we can learn a good enough approximation to the histogram using $\tilde{O}\left(\sqrt{N}\right) \mathsf{poly}(\tau^{-1})$ samples from $D$.

For $X_{light}$, if our estimate on its probability is less than least $\sigma/9$, then we can safely ignore it. Otherwise, similarly to the above, we know that w.h.p. $X_{light}$'s probability by $D$ is at least $\sigma/10$. In this case, we run a verified histogram protocol outlined in the previous section to obtain the tagged sample $(z_i, \pi(z_i))_{i \in [s]}$ of the distribution $D\big|_{X_{light}}$, from which we can deduce the histogram of $D\big|_{X_{light}}$, up to $\Delta_{\mathrm{RL}}$ distance $\sigma/10$. Observe that we can use rejection sampling to implement sample-access to $D\big|_{X_{light}}$ via samples form $D$ (with an overhead of at most $O(1/\sigma)$ samples from $D$ per sample from $D\big|_{X_{light}}$), so we can indeed run the protocol on $D\big|_{X_{light}}$ (paying the aforementioned overhead).

To wrap up, if our estimate on the weight of either of the two sets is smaller than $\sigma/9$, then w.h.p. its true weight is less than $\sigma/8$, and we can safely ignore it. We learn the other one to within distance $\sigma/10$, so the total distance from the learned histogram to the true one is smaller than $0.9\sigma$ with high probability. If both sets have empirical weight greater than $\sigma/10$, then we learn the two histograms and compose them into a single global histogram. W.h.p. the result will

---

[4]In fact, we learn the explicit distribution $D'$ over domain $X_{heavy}$, guaranteed to be $\sigma/10$ close in statistical distance to $D\big|_{X_{heavy}}$ with high probability, not just the histogram

be at distance smaller than $0.9\sigma$ from the true distribution of $D'$. For more details, we refer the reader to the fuller exposition in [HR22], and the *IP for verified histogram reconstruction.*

## 6.3    From Doubly Efficient Histogram Verification to Doubly Efficient Verification of Label-Invariant Distribution Properties

In order to get from doubly efficient histogram reconstruction protocol, as presented here, to a proof system for tolerantly-verifying a large collection of label invariant distribution properties, the reader is referred to Section 4 of [HR22]. We explain here how this is achieved, but omit details.

Recall the setting of Theorem 1.1. Let $\mathcal{P}$ be some label-invariant distribution property. The verifier get as input $N \in \mathbb{N}$, parameters $0 \leq \varepsilon_c < \varepsilon_f \leq 1$, as well as black-box sample access to distribution $D$ over domain $[N]$. Denote $\rho = \varepsilon_f - \varepsilon_c$. We want a proof system that with high probability achieves the following:

- **Completeness.** If $\Delta_{\mathrm{SD}}(D, \mathcal{P}) \leq \varepsilon_c$, the verifier accepts with high probability.

- **Soundness.** If $\Delta_{\mathrm{SD}}(D, \mathcal{P}) \geq \varepsilon_f$, the verifier rejects with high probability.

- The verifier sample complexity, runtime, and the protocol communication complexity are all $\widetilde{O}\left(\sqrt{N}\right) \mathsf{poly}\left(\rho^{-1}\right)$.

- The prover runs in time $\widetilde{O}(N)\mathsf{poly}\left(\rho^{-1}\right)$

Consider the following protocol:

1. **Obtain a verified approximate histogram of $D$.** Run the protocol outlined in Section 6.2 with $\sigma = \frac{\rho}{3}$, and $\tau = \frac{1}{500}\sigma^2$. At the end of the run, either the verifier rejected, or they have obtained an $(N, \tau)$- histogram $\{a_j\}_j$ that satisfies $\Delta_{\mathrm{RL}}(D, \{a_j\}_j) \leq \rho/3$.

The verifier now knows a good approximation of the histogram of $D$. From here, it needs to determine whether $D$ is $\varepsilon_c$-close or $\varepsilon_f$-far from $\mathcal{P}$. This is done through the following procedure:

2. **Obtain a histogram of a distribution inside $\mathcal{P}$ close to $D$.** The prover provides a histogram $\{b_j\}_j$ so that there exists a distribution $D'$ consistent with $\{b_j\}_j$ such that $D' \in \mathcal{P}$, and $\Delta_{\mathrm{SD}}(D', D) \leq \varepsilon_c$.

Now, the verifier needs to figure out two things: $(i)$ whether $\{b_j\}_j$ is indeed consistent with some distribution inside $\mathcal{P}_N$. If so, since $\mathcal{P}$ is label-invariant, it means that *all* distributions consistent with $\{b_j\}_j$ are in $\mathcal{P}_N$ (or $\tau$-close to one); $(ii)$ whether there exists a distribution $D'$ consistent with $\{b_j\}_j$ which is $\varepsilon_c$-close to $D$.

The second point can be easily checked through an algorithm given in [HR22]. The first point however is tricky. In order to efficiently produce $\{b_j\}_j$ by the prover and in order to find a procedure that checks that $\{b_j\}_j$ is consistent with some distribution in $\mathcal{P}$ we actually limit ourselves to only those label-invariant distribution properties that admit such procedures. Formally:

**Definition 6.5** (Doubly-efficient approximate decision procedure)**.** *For every $N \in \mathbb{N}$, and distribution property $\mathcal{P}$, denote $\mathcal{P}_N = \mathcal{P} \cap \Delta_N$. $\mathcal{P}$ has* doubly-efficient approximate decision procedures *if there exist two algorithms $A_{check}$ and $A_{produce}$ as follows:*

Procedure $A_{check}$ runs in time $\mathsf{poly}(\log N, \tau^{-1})$, gets as input the domain size $N$, accuracy parameter $\tau$, distance parameter $\sigma \in (0,1)$, and a $(N,\tau)$-histogram $\{p_j\}_j$. There exists a function $\mu(N,\sigma) = \mathsf{poly}(1/\log N, \sigma)$ ($\mu$ specifies a minimum sensitivity requirement on the accuracy parameter $\tau$) s.t. for every integer $N$, $\tau \leq \mu(N,\sigma)$, $(N,\tau)$-histogram $\{p_j\}_j$, and $\sigma > 0$:

- If there exists a distribution $D \in \mathcal{P}_N$ consistent with $\{p_j\}_j$, then $A_{check}$ accepts.

- If every distribution $D \in \Delta_N$ consistent with $\{p_j\}_j$ satisfies $\Delta_{RL}(D, \mathcal{P}_N) \geq \sigma$, $A_{check}$ rejects.

$A_{produce}$ gets black-box sample access to the distribution $D$, which is assumed to satisfy $\Delta_{SD}(D, \mathcal{P}_N) \leq \varepsilon_c$, as well as parameters $N, \tau, \varepsilon_c$, and $\varepsilon_f$, it has runtime and sample complexity of magnitude $\widetilde{O}(N)\mathsf{poly}\left((\varepsilon_f - \varepsilon_c)^{-1}\right)$, and produces a histogram $\{b_j\}_j$ such that: $\{b_j\}_j$ is consistent with some $D'' \in \mathcal{P}_N$ and $\Delta_{RL}(D, \{b_j\}_j) \leq \varepsilon_c + \frac{\varepsilon_f - \varepsilon_c}{10}$.

And so, if we assume $\mathcal{P}$ admits *doubly-efficient approximate decision procedures*, the honest prover can produce the histogram $\{b_j\}_j$ efficiently, and the verifier can check that $\{b_j\}_j$ is close to $\{a_j\}_j$, run the decision procedure, and accept if both checks pass.

**Completeness and Soundness.** Completeness is immediate from the completeness of all steps involved. Consider next that $D$ is $\varepsilon_f$-far from $\mathcal{P}$. Then, if the verifier in the protocol outlined in Section 6.2 didn't reject the histogram obtained, then with high probability $\Delta_{\mathrm{RL}}(D, \{a_j\}_j) \leq \sigma = \rho/3$. Then, consider $\{b_j\}_j$, if there doesn't exist a distribution consistent with $\{b_j\}_j$ inside $\mathcal{P}_N$, the verifier will reject. Therefore, assume that there is such a distribution. Next, it must hold that $\Delta_{\mathrm{RL}}(\{b_j\}_j, D) \geq \varepsilon_f$, which means that: $\Delta_{\mathrm{RL}}(\{a_j\}_j, \{b_j\}_j) \geq \Delta_{\mathrm{RL}}(\{b_j\}_j, D) - \Delta_{\mathrm{RL}}(\{a_j\}_j, D) \geq \varepsilon_f - \rho > \varepsilon_c + \rho/3$. And we get that the test fails. And the verifier rejects with high probability

**Prover Complexity.** The prover has to run two procedures: 1. Protocol 6.1.1, as part of the proof system outlined in Section 6.2. This requires it to run in time $\widetilde{O}(N)\mathsf{poly}(\rho^{-1})$; 2. It needs to provide $\{b_j\}_j$, which can be produced efficiently through $A_{produce}$ in time $\widetilde{O}(N)\mathsf{poly}(\rho^{-1})$, assuming that $\mathcal{P}$ has a doubly-efficient approximate decision procedure.

We remark that producing proofs for properties with $A_{produce}$ that runs in time super-linear is possible as well, but it will result with higher runtime for the prover.

**About doubly-efficient decision procedures.** We consider the existence of *doubly efficient decision procedures* to be a mild assumption on the distribution property. The reader is referred to [HR22] for more further details on the check-procedure, as well as examples of explicit algorithms for natural distribution properties.

In this work, we also require the existence of $A_{produce}$, which was not required in [HR22]. We consider this additional requirement as an only slightly stronger assumption. This is since in time quasi-linear in $N$, the honest prover can obtain an explicit description of a distribution $D'$ that is $\frac{\varepsilon_f - \varepsilon_c}{100}$-close to $D$. A possible implementation of $A_{produce}$ can simply change $D'$ to obtain some description of $D''$ (the histogram of which will be sent to the verifier) such that $\Delta_{\mathrm{SD}}(D'', D') \leq \varepsilon_c + \frac{\varepsilon_f - \varepsilon_c}{100}$, and $D'' \in \mathcal{P}_N$. For properties such as being at distance $\delta$ from $U_{[N]}$, or having entropy smaller than $K$, the process of obtaining $D''$ for $D'$ requires only $\widetilde{O}(N)$ runtime (indeed, for these properties, a histogram of $D''$ can even be obtained in *sublinear* time by computing it directly form the histogram of $D'$).

# 7 Acknowledgment

We thank Oded Goldreich for illuminating conversations that were tremendously helpful in improving our understanding and the presentation of this work

# References

[BC17]    Tugkan Batu and Clément L. Canonne. Generalized uniformity testing. In Chris Umans, editor, *58th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2017, Berkeley, CA, USA, October 15-17, 2017*, pages 880–889. IEEE Computer Society, 2017.

[BFF+01]  Tugkan Batu, Lance Fortnow, Eldar Fischer, Ravi Kumar, Ronitt Rubinfeld, and Patrick White. Testing random variables for independence and identity. In *42nd Annual Symposium on Foundations of Computer Science, FOCS 2001, 14-17 October 2001, Las Vegas, Nevada, USA*, pages 442–451. IEEE Computer Society, 2001.

[BFLS91]  László Babai, Lance Fortnow, Leonid A. Levin, and Mario Szegedy. Checking computations in polylogarithmic time. In Cris Koutsougeras and Jeffrey Scott Vitter, editors, *Proceedings of the 23rd Annual ACM Symposium on Theory of Computing, May 5-8, 1991, New Orleans, Louisiana, USA*, pages 21–31. ACM, 1991.

[CG18]    Alessandro Chiesa and Tom Gur. Proofs of proximity for distribution testing. In Anna R. Karlin, editor, *9th Innovations in Theoretical Computer Science Conference, ITCS 2018, January 11-14, 2018, Cambridge, MA, USA*, volume 94 of *LIPIcs*, pages 53:1–53:14. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2018.

[EKR04]   Funda Ergün, Ravi Kumar, and Ronitt Rubinfeld. Fast approximate probabilistically checkable proofs. *Inf. Comput.*, 189(2):135–159, 2004.

[GGR98]   Oded Goldreich, Shafi Goldwasser, and Dana Ron. Property testing and its connection to learning and approximation. *J. ACM*, 45(4):653–750, 1998.

[GKR15]   Shafi Goldwasser, Yael Tauman Kalai, and Guy N. Rothblum. Delegating computation: Interactive proofs for muggles. *J. ACM*, 62(4):27:1–27:64, 2015.

[GMR85]   Shafi Goldwasser, Silvio Micali, and Charles Rackoff. The knowledge complexity of interactive proof-systems (extended abstract). In Robert Sedgewick, editor, *Proceedings of the 17th Annual ACM Symposium on Theory of Computing, May 6-8, 1985, Providence, Rhode Island, USA*, pages 291–304. ACM, 1985.

[GMW91]   Oded Goldreich, Silvio Micali, and Avi Wigderson. Proofs that yield nothing but their validity for all languages in NP have zero-knowledge proof systems. *J. ACM*, 38(3):691–729, 1991.

[Gol17]   Oded Goldreich. *Introduction to Property Testing*. Cambridge University Press, 2017.

[Gol20]   Oded Goldreich. The uniform distribution is complete with respect to testing identity to a fixed distribution. In Oded Goldreich, editor, *Computational Complexity and Property*

Testing - On the Interplay Between Randomness and Computation, volume 12050 of *Lecture Notes in Computer Science*, pages 152–172. Springer, 2020.

[GR18]      Tom Gur and Ron D. Rothblum.  Non-interactive proofs of proximity.  *Comput. Complex.*, 27(1):99–207, 2018.

[GRSY21]  Shafi Goldwasser, Guy N. Rothblum, Jonathan Shafer, and Amir Yehudayoff. Interactive proofs for verifying machine learning.  In James R. Lee, editor, *12th Innovations in Theoretical Computer Science Conference, ITCS 2021, January 6-8, 2021, Virtual Conference*, volume 185 of *LIPIcs*, pages 41:1–41:19. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2021.

[HILL99]   Johan Håstad, Russell Impagliazzo, Leonid A. Levin, and Michael Luby.  A pseudorandom generator from any one-way function. *SIAM J. Comput.*, 28(4):1364–1396, 1999.

[HR22]      Tal Herman and Guy N. Rothblum. Verifying the unseen: interactive proofs for label-invariant distribution properties.  In Stefano Leonardi and Anupam Gupta, editors, *STOC '22: 54th Annual ACM SIGACT Symposium on Theory of Computing, Rome, Italy, June 20 - 24, 2022*, pages 1208–1219. ACM, 2022.

[Mic94]     Silvio Micali.  CS proofs (extended abstracts).  In *35th Annual Symposium on Foundations of Computer Science, Santa Fe, New Mexico, USA, 20-22 November 1994*, pages 436–453. IEEE Computer Society, 1994.

[PRR06]    Michal Parnas, Dana Ron, and Ronitt Rubinfeld. Tolerant property testing and distance approximation. *J. Comput. Syst. Sci.*, 72(6):1012–1042, 2006.

[RS96]      Ronitt Rubinfeld and Madhu Sudan.  Robust characterizations of polynomials with applications to program testing. *SIAM J. Comput.*, 25(2):252–271, 1996.

[RVW13]   Guy N. Rothblum, Salil P. Vadhan, and Avi Wigderson. Interactive proofs of proximity: delegating computation in sublinear time. In Dan Boneh, Tim Roughgarden, and Joan Feigenbaum, editors, *Symposium on Theory of Computing Conference, STOC'13, Palo Alto, CA, USA, June 1-4, 2013*, pages 793–802. ACM, 2013.

[SV03]       Amit Sahai and Salil P. Vadhan. A complete problem for statistical zero knowledge. *J. ACM*, 50(2):196–249, 2003.

[Vad99]     Salil Vadhan. *A Study of Statistical Zero-Knowledge Proofs*. PhD thesis, USA, 1999.

[Vad12]     Salil Vadhan. *Pseudorandomness*. Now Publishers Inc., Hanover, MA, USA, 2012.

[VV10]       Gregory Valiant and Paul Valiant. A CLT and tight lower bounds for estimating entropy. *Electron. Colloquium Comput. Complex.*, 17:183, 2010.

[VV14]       Gregory Valiant and Paul Valiant. An automatic inequality prover and instance optimal identity testing. In *55th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2014, Philadelphia, PA, USA, October 18-21, 2014*, pages 51–60. IEEE Computer Society, 2014.

# A    Collisions Concentration

Assume $D$ is a distribution over domain $[N]$ that satisfies that for every $x \in [N]$, $D(x) \leq \frac{\tau}{\sqrt{N}}$ with parameters. Let $S$ be an i.i.d. sample of size $s$ drawn by $D$. Let $\{\widehat{q}_\ell\}_\ell, \{x_{\ell,j}\}_{\ell,j}$ be as defined in Section 4. Let $J$ be as in Definition 5.2

Assume a fresh i.i.d. sample of $D$ of size $s$ was sampled. Denote this sample by $T$. And for every bucket $j$, define $\widetilde{C}_j$ to be defined in Tester 5.1.1.

This section is taken with small changes from [HR22]. The main difference is that the notation is according to the conventions in this paper, and also, we consider only buckets satisfying the conditions of set $J$.

**Claim A.1.** *With probability of at least* 0.99 *over the choice of sample $S$, and any mislabelling of $S$ characterised by variables $\{x_{\ell,j}\}_{\ell,j}$:*

- *For every $j \in J$:*

$$
\mathbb{E}[\widetilde{C}_j] \in \frac{s^2}{N} \left[ \sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{\ell\tau}, e^{2\tau} \sum_{\ell:e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{\ell\tau} \right)
$$

- *With probability of at least 0.99 over the choice of $T$, for all $j$ in $J$:*

$$
\left| \widetilde{C}_j - \mathbb{E}_T[\widetilde{C}_j] \right| \leq \tau \mathbb{E}_T[\widetilde{C}_j]
$$

*Proof.* Fix $j \in J$. Define $I_{r,k}$ to be the indicator that $T_r = S_k$, and $\text{tag}(S_k) = j$. Denote $\widetilde{F}_j = \{i \in [s] : \text{tag}(S_i) = j\}$. By definition, $\widetilde{C}_j = \sum_{r \in [s]} \sum_{k \in \widetilde{F}_j} I_{r,k}$ (note that $I_{r,k} = 0$ for all $k \notin \widetilde{F}_j$). Therefore, by the linearity of expectation:

$$
\mathbb{E}\left[\widetilde{C}_j\right] = \sum_{r \in [s]} \sum_{k \in \widetilde{F}_j} \mathbb{E}[I_{r,k}]
$$

The value of $\mathbb{E}[I_{r,k}]$ can vary significantly between indices in $\widetilde{F}_j$, depending on $S_k$, the probability of the element $S_k$ affects the probability that the sample $T_r$ collided with it. Thus, we divide $\widetilde{F}_j$ into disjoint subsets according to the bucket origin of each sample in $S$. Define $\widetilde{F}_{\ell \rightarrow j} \subseteq \widetilde{F}_j$ to be the set of indices associated with true bucket $\ell$ that were tagged as belonging to alleged bucket $j$. By this definition, $\widetilde{F}_j = \cup_\ell \widetilde{F}_{\ell \rightarrow j}$, and also for every $\ell$, $|\widetilde{F}_{\ell \rightarrow j}| = s\widehat{q}_\ell x_{\ell,j}$. Plugging this back to the above expression:

$$
\mathbb{E}[\widetilde{C}_j] = \sum_\ell \sum_{r \in [s]} \sum_{k \in \widetilde{F}_{\ell \rightarrow j}} \mathbb{E}[I_{r,k}] = \sum_{r \in [s]} \sum_{k \in \widetilde{F}_{L \rightarrow j}} \mathbb{E}[I_{r,k}] + \sum_{\ell:e^{\ell\tau} \geq \tau^2} \sum_{r \in [s]} \sum_{k \in \widetilde{F}_{\ell \rightarrow j}} \mathbb{E}[I_{r,k}] \tag{50}
$$

This decomposition of the sum allows us to unravel the expression $\mathbb{E}[I_{r,k}]$, since for all $\ell$, that satisfy $e^{\ell\tau} \geq \tau^2$, every $k \in \widetilde{F}_{\ell \rightarrow j}$, satisfies $D(S_k) \in \left[ \frac{e^{\ell\tau}}{N}, e^\tau \frac{e^{\ell\tau}}{N} \right)$, and so $\mathbb{E}[I_{r,k}] \in \left[ \frac{e^{\ell\tau}}{N}, e^\tau \frac{e^{\ell\tau}}{N} \right)$. Similarly, for $k \in \widetilde{F}_{L \rightarrow j}$, $D(S_k) = \mathbb{E}[I_{r,k}] \in \left[ 0, \frac{\tau^2}{N} \right)$. We conclude that:

$$
\sum_{\ell:e^{\ell\tau} \geq \tau^2} \sum_{r \in [s]} \sum_{k \in \widetilde{F}_{\ell \rightarrow j}} \mathbb{E}[I_{r,k}] \leq \sum_{\ell:e^{\ell\tau} \geq \tau^2} s \cdot (sx_{\ell,j}\widehat{q}_\ell) \cdot e^\tau \frac{e^{\ell\tau}}{N} = e^\tau \sum_{\ell:e^{\ell\tau} \geq \tau^2} x_{\ell,j} \cdot \frac{s^2}{N} \widehat{q}_\ell e^{\ell\tau} \tag{51}
$$

And:

$$\sum_{\ell:e^{\ell\tau}\geq\tau^2}\sum_{r\in[s]}\sum_{k\in\widetilde{F}_{\ell\to j}}\mathbb{E}[I_{r,k}]\geq\sum_{\ell:e^{\ell\tau}\geq\tau^2}s\cdot(sx_{\ell,j}\widehat{q}_\ell)\cdot\frac{e^{\ell\tau}}{N}=\sum_{\ell:e^{\ell\tau}\geq\tau^2}x_{\ell,j}\cdot\frac{s^2}{N}\widehat{q}_\ell e^{\ell\tau} \tag{52}$$

As well as:

$$\sum_{r\in[s]}\sum_{k\in\widetilde{F}_{L\to j}}\mathbb{E}[I_{r,k}]\leq s\cdot(sx_{L,j}\widehat{p}_L)\cdot\frac{\tau^2}{N}=\frac{s^2}{N}\widehat{q}_L x_{L,j}\tau^2 \tag{53}$$

Finally, observe that given $j\in J$, it holds that $\widehat{q}_L\leq\tau v_j$, which implies:

$$\widehat{q}_L x_{L,j}\leq\tau v_j=\tau\sum_{\ell:e^{\ell\tau}\geq\tau^2}\widehat{q}_\ell x_{\ell,j}+\tau\cdot\widehat{q}_L$$

Or equivalently, since $\tau<0.01$:

$$\widehat{q}_L x_{L,j}\leq\frac{\tau}{1-\tau}\sum_{\ell:e^{\ell\tau}\geq\tau^2}\widehat{q}_\ell x_{\ell,j}\leq1.1\tau\sum_{\ell:e^{\ell\tau}\geq\tau^2}\widehat{q}_\ell x_{\ell,j}$$

From which we conclude:

$$\frac{s^2}{N}\widehat{q}_L x_{L,j}\tau^2\leq1.1\tau\sum_{\ell:e^{\ell\tau}\geq\tau^2}\widehat{q}_\ell x_{\ell,j}\tau^2\leq1.1\tau\sum_{\ell:e^{\ell\tau}\geq\tau^2}\widehat{q}_\ell x_{\ell,j}e^{\ell\tau} \tag{54}$$

Combining Inequalities (53),(52), and (54), we conclude:

$$\mathbb{E}[\widetilde{C}_j]\leq e^\tau\sum_{\ell:e^{\ell\tau}\geq\tau^2}\left(\frac{s^2}{N}\widehat{q}_\ell x_{\ell,j}e^{\ell\tau}\right)+\frac{s^2}{N}\widehat{q}_L x_{L,j}\tau^2\leq e^{2\tau}\sum_{\ell:e^{\ell\tau}\geq\tau^2}\left(\frac{s^2}{N}\widehat{q}_\ell x_{\ell,j}e^{\ell\tau}\right)$$

We also get:

$$\mathbb{E}[\widetilde{C}_j]\geq\sum_{\ell:e^{\ell\tau}\geq\tau^2}\frac{s^2}{N}\widehat{q}_\ell x_{\ell,j}e^{\ell\tau}$$

And so concludes the first part of the proof.

Moving on to proving measure concentration. In order to do so, we bound $\mathrm{Var}[\widetilde{C}_j]$ from above, in the aim of using Chebyshev's inequality to bound the probability that $\widetilde{C}_j$ deviates from its expectation. First, recall that:

$$\mathrm{Var}\left[\widetilde{C}_j\right]=\sum_{\substack{(r_1,k_1):\\r_1\in[s]\\k_1\in\widetilde{F}_j}}\sum_{\substack{(r_2,k_2):\\r_2\in[s]\\k_2\in\widetilde{F}_j}}\mathrm{Cov}\left[I_{r_1,k_1},I_{I_2,k_2}\right]$$

In order to bound this expression, observe that for every $r_1,r_2\in[s]$, such that $r_1\neq r_2$, since $T_{r_1}$ and $T_{r_2}$ were chosen i.i.d., the variables $I_{r_1,k_1}$ and $I_{r_2,k_2}$ are independent, and so $\mathrm{Cov}\left[I_{r_1,k_1},I_{r_2,k_2}\right]=0$. Also, if $r_1=r_2$, but $S_{k_1}\neq S_{k_2}$, then, as it is impossible that both the variables $I_{r_1,k_1},I_{r_2,k_2}$ are positive at the same time, it follows that in this case, $\mathrm{Cov}\left[I_{r_1,k_1},I_{r_2,k_2}\right]<0$. This leaves us

only with the case $r_1 = r_2$ and $S_{k_1} = S_{k_2}$. In this case, the variables satisfy $I_{r_1,k_1} = I_{r_2,k_2}$, and by the definition of the covariance, this yields $\text{Cov}[I_{r_1,k_1}, I_{r_2,k_2}] = \text{Var}[I_{r_1,k_1}]$. And as for every $r_1, k_1$, $\text{Var}[I_{r_1,k_1}] \leq \mathbb{E}_T[I_{r_1,k_1}]$, we conclude:

$$\text{Var}[\widetilde{C}_j] \leq \sum_{r \in [s]} \sum_{k_1 \in \widetilde{F}_j} \sum_{\substack{k_2: \\ S_{k_2} = S_{k_1}}} \mathbb{E}_T[I_{r,k_1}] \tag{55}$$

Assuming $D$ has maximal probability $\tau/\sqrt{N}$, with probability at least $0.99$ over the choice of $S$, it follows that every element sampled in $S$ appears at most $\log N$ times, and so, we are guaranteed that for every $k_1$, the number of summands in the third sum over $k_2$ is at most $\log N$. Therefore:

$$\text{Var}[\widetilde{C}_j] \leq \sum_{r \in [s]} \sum_{k_1 \in \widetilde{F}_\ell} \sum_{\substack{k_2: \\ S_{k_2} = S_{k_1}}} \mathbb{E}_T[I_{r,k_1}] \leq \log N \sum_{r \in [s]} \sum_{k_1 \in \widetilde{F}_j} \mathbb{E}_T[I_{r,k_1}] = \log N \mathbb{E}_T[\widetilde{C}_j] \tag{56}$$

For every $j \in J$:

$$\mathbb{E}\left[\widetilde{C}_j\right] \geq \frac{s^2}{N} \sum_{\ell: e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} e^{\ell\tau} \geq \frac{s^2 \cdot \tau^2}{N} \frac{s^2}{N} \sum_{\ell: e^{\ell\tau} \geq \tau^2} \widehat{q}_\ell x_{\ell,j} \geq \frac{s^2 \cdot \tau^2}{N} v_j(1-\tau) \geq \frac{s^2}{N} \cdot \frac{\tau^4}{\log N}$$

Therefore, using Chebyshev's inequality, as well as the lower bound for shown above, for every $j \in J$:

$$\Pr_T\left(\left|\widetilde{C}_j - \mathbb{E}_T[\widetilde{C}_j]\right| \geq \tau \mathbb{E}_T[\widetilde{C}_j]\right) \leq \frac{\log N}{\tau^2 \mathbb{E}_T[\widetilde{C}_j]} \leq \frac{N \log^2 N}{\tau^6} \cdot \frac{1}{s}$$

Summing over all $j \in J$, we get by union bound that the probability that there exists some $j \in J$ such that $\left|\widetilde{C}_\ell - \mathbb{E}_T[\widetilde{C}_\ell]\right| > \tau \mathbb{E}_T[\widetilde{C}_\ell]$ is at most:

$$\sum_{j \in J} \frac{N \log^2 N}{\tau^6} \cdot \frac{1}{s} \leq \frac{b(N,\tau) \cdot N \log^2 N}{\tau^6} \cdot \frac{1}{s} \leq \frac{1}{s} \leq 0.01$$

Where the last inequality is justified by the choice of $s$. $\qquad\square$

**Corollary A.2.** *If $x_{\ell,j} = \mathbb{1}_{\ell=j}$ for all $(\ell, j)$ then, for all $j$ such that $e^{j\tau} \geq \tau$ and $v_j \geq \frac{\tau^2}{\log N}$ it holds that with probability at least $0.99$ over the choice of $S$ and $T$:*

$$\left|\widetilde{C}_j - \frac{s^2}{N} v_j e^{j\tau}\right| \leq 5\tau \frac{s^2}{N} v_j e^{j\tau}$$

*Proof.* Plugging $x_{\ell,j} = \mathbb{1}_{\ell=j}$ in Claim A.1, we get that every $j$ for which $e^{j\tau} \geq \tau$ and $v_j \geq \frac{\tau^2}{\log N}$ also satisfies $\widehat{q}_L x_{L,j} = 0 \leq \tau v_j$, and so $j \in J$. Therefore, with probability at least $0.99$, for all such

$j$, since $\tau < 0.01$:

$$\left| \widetilde{C}_j - \frac{s^2}{N} v_j e^{j\tau} \right| \leq \left| \widetilde{C}_j - \mathbb{E}_T[\widetilde{C}_j] \right| + \left| \mathbb{E}_T[\widetilde{C}_j] - \frac{s^2}{N} v_j e^{j\tau} \right|$$

$$\leq \tau \mathbb{E}_T[\widetilde{C}_j] + \left| \mathbb{E}_T[\widetilde{C}_j] - \frac{s^2}{N} v_j e^{j\tau} \right|$$

$$\leq \tau \cdot e^{2\tau} \frac{s^2}{N} v_j e^{j\tau} + \left( e^{2\tau} - 1 \right) \frac{s^2}{N} v_j e^{j\tau}$$

$$\leq \left( \tau e^{2\tau} + e^{2\tau} - 1 \right) \frac{s^2}{N} v_j e^{j\tau}$$

$$\leq \left( \tau + 3\tau^2 + 3\tau \right) \frac{s^2}{N} v_j e^{j\tau}$$

$$\leq 5\tau \frac{s^2}{N} v_j e^{j\tau}$$

$\square$