



Testing probability distributions using conditional samples

Clément Canonne* Dana Ron† Rocco A. Servedio‡

November 12, 2012

Abstract

We study a new framework for property testing of probability distributions, by considering distribution testing algorithms that have access to a *conditional sampling oracle*.¹ This is an oracle that takes as input a subset $S \subseteq [N]$ of the domain $[N]$ of the unknown probability distribution D and returns a draw from the conditional probability distribution D restricted to S . This new model allows considerable flexibility in the design of distribution testing algorithms; in particular, testing algorithms in this model can be adaptive.

We study a wide range of natural distribution testing problems in this new framework and some of its variants, giving both upper and lower bounds on query complexity. These problems include testing whether D is the uniform distribution \mathcal{U} ; testing whether $D = D^*$ for an explicitly provided D^* ; testing whether two unknown distributions D_1 and D_2 are equivalent; and estimating the variation distance between D and the uniform distribution. At a high level our main finding is that the new “conditional sampling” framework we consider is a powerful one: while all the problems mentioned above have $\Omega(\sqrt{N})$ sample complexity in the standard model (and in some cases the complexity must be almost linear in N), we give $\text{poly}(\log N, 1/\epsilon)$ -query algorithms (and in some cases $\text{poly}(1/\epsilon)$ -query algorithms independent of N) for all these problems in our conditional sampling setting.

*cannon@cs.columbia.edu, Columbia University.

†danaron@post.tau.ac.il, Tel Aviv University. Supported by ISF grant number 246/08.

‡rocco@cs.columbia.edu, Columbia University. Supported by NSF grants CCF-0915929 and CCF-1115703.

¹Independently from our work, Chakraborty et al. [CFG13] also considered this framework. We discuss their work in Subsection 1.4.

1 Introduction

1.1 Background: Distribution testing in the standard model

One of the most fundamental problem paradigms in statistics is that of inferring some information about an unknown probability distribution D given access to independent samples drawn from it. More than a decade ago, Batu et al. [BFR⁺00]² initiated the study of problems of this type from within the framework of *property testing* [RS96, GGR98]. In a property testing problem there is an unknown “massive object” that an algorithm can access only by making a small number of “local inspections” of the object, and the goal is to determine whether the algorithm has a particular property. The algorithm must output ACCEPT if the object has the desired property and output REJECT if the object is far from every object with the property. (See [Fis01, Ron08, Ron10, Gol10] for detailed surveys and overviews of the broad field of property testing; we give precise definitions tailored to our setting in Section 2.)

In distribution property testing the “massive object” is an unknown probability distribution D over an N -element set, and the algorithm accesses the distribution by drawing independent samples from it. A wide range of different properties of probability distributions have been investigated in this setting, and upper and lower bounds on the number of samples required have by now been obtained for many problems. These include testing whether D is uniform [GR00, BFR⁺10, Pan08], testing whether D is identical to a given known distribution D^* [BFF⁺01], testing whether two distributions D_1, D_2 (both available via sample access) are identical [BFR⁺00, Val11], and testing whether D has a monotonically increasing probability mass function [BFRV11], as well as related problems such as estimating the entropy of D [BDKR05, VV11], and estimating its support size [RRSS09, Val11, VV11]. Similar problems have also been studied by researchers in other communities, see e.g., [Ma81, Pan04, Pan08].

One broad insight that has emerged from this past decade of work is that while sublinear-sample algorithms do exist for many distribution testing problems, the number of samples required is in general quite large. Even the basic problem of testing whether D is the uniform distribution \mathcal{U} over $[N] = \{1, \dots, N\}$ versus ϵ -far from uniform requires $\Omega(\sqrt{N})$ samples³ for constant ϵ , and the other problems mentioned above have sample complexities at least this high, and in some cases *almost linear in N* [RRSS09, Val11, VV11]. Since such sample complexities could be prohibitively high in real-world settings where N can be extremely large, it is natural to explore problem variants where it may be possible for algorithms to succeed using fewer samples. Indeed, researchers have studied distribution testing in settings where the unknown distribution is guaranteed to have some special structure, such as being monotone, k -modal or a “ k -histogram” over $[N]$ [BKR04, DDS⁺13, ILR12], or being monotone over $\{0, 1\}^n$ [RS09] or over other posets [BFRV11], and have obtained significantly more sample-efficient algorithms using these additional assumptions.

²There is a more recent full version of this work [BFR⁺10] and we henceforth reference this recent version.

³To verify this, consider the family of all distributions that are uniform over half of the domain, and 0 elsewhere. Each distribution in this family is $\Theta(1)$ -far from the uniform distribution. However, it is not possible to distinguish with sufficiently high probability between the uniform distribution and a distribution selected randomly from this family, given a sample of size \sqrt{N}/c (for a sufficiently large constant $c > 1$). This is the case because for the uniform distribution as well as each distribution in this family, the probability of observing the same element more than once is very small. Conditioned on such a collision event not occurring, the samples are distributed identically.

1.2 Our model: Conditional sampling

In this work we pursue a different line of investigation: rather than restricting the class of probability distributions under consideration, we consider testing algorithms that may use a more powerful form of access to the unknown distribution D . This is a *conditional sampling oracle*, which allows the algorithm to obtain a draw from D_S , the conditional distribution of D restricted to a subset S of the domain, (where S is specified by the algorithm). More precisely, we have:

Definition 1 *Fix a distribution D over $[N]$. A COND oracle for D , denoted COND_D , is defined as follows: The oracle is given as input a query set $S \subseteq [N]$ that has $D(S) > 0$. The oracle returns an element $i \in S$, where the probability that element i is returned is $D_S(i) = D(i)/D(S)$, independently of all previous calls to the oracle.⁴*

For compatibility with our COND_D notation we will write SAMP_D to denote an oracle that takes no input and, each time it is invoked, returns an element from $[N]$ drawn according to D independently from all previous draws. This is the sample access to D that is used in the standard model of testing distributions, and this is of course the same as a call to $\text{COND}_D([N])$.

Motivation and Discussion. The first motivation for studying the COND model is to capture scenarios that arise in application areas (e.g., in biology or chemistry) where the parameters of some experiment can be adjusted so as to restrict the range of possible outcomes. For example, a scientist growing bacteria or yeast cells in a controlled environment may be able to deliberately introduce environmental factors that allow only cells with certain desired characteristics to survive, corresponding to restricting the distribution of all experimental outcomes to a pre-specified subset.

A second, purely theoretical motivation, is that the study of the COND model may further our understanding regarding what forms of information (beyond standard sampling) can be helpful for testing properties of distributions. In both learning and property testing it is generally interesting to understand how much power algorithms can gain by making queries, and COND queries are a natural type of query to investigate in the context of distributions. As we discuss in more detail below, in several of our results we actually consider restricted versions of COND queries that do not require the full power of obtaining conditional samples from arbitrary sets.

A third attractive feature of the COND model is that it enables a new level of “richness” for algorithms that deal with probability distributions. In the standard model where only access to SAMP_D is provided, all algorithms must necessarily be non-adaptive, with the same initial step of simply drawing a sample of points from SAMP_D , and the difference between two algorithms comes only from how they process their samples. In contrast, the essence of the COND model is to allow algorithms to *adaptively* determine later query sets S based on the outcomes of earlier queries.

⁴Note that as described above the behavior of $\text{COND}_D(S)$ is undefined if $D(S) = 0$, i.e., the set S has zero probability under D . While various definitional choices could be made to deal with this, we shall assume that in such a case, the oracle (and hence the algorithm) outputs “failure” and terminates. This will never be a problem for us throughout this paper, as (a) all of our lower bounds will deal only with distributions that have $D(i) > 0$ for all $i \in [N]$, and (b) in all of our algorithms $\text{COND}_D(S)$ will only ever be called on sets S which are “guaranteed” to have $D(S) > 0$. (More precisely, each time an algorithm calls $\text{COND}_D(S)$ it will either be on the set $S = [N]$, or will be on a set S which contains an element i which has been returned as the output of an earlier call to COND_D .)

Given the above motivations, the central question is whether the COND model enables significantly more efficient algorithms than are possible in the weaker SAMP model. Our results (see Section 1.3) show that this is indeed the case.

Before detailing our results, we note that many of our results will in fact deal with a weaker variant of the COND model, which we now describe. In designing COND-model algorithms it is obviously desirable to have algorithms that only invoke the COND oracle on query sets S which are “simple” in some sense. Of course there are many possible notions of “simplicity”; in this work we consider two such notions, corresponding to two restrictions of the general COND model, as follows:

- **PCOND oracle:** We define a PCOND (short for “pair-cond”) oracle for D as a restricted version of COND_D that only accepts input sets S which are either $S = [N]$ (thus providing the power of a SAMP_D oracle) or $S = \{i, j\}$ for some $i, j \in [N]$, i.e. sets of size two. The PCOND oracle may be viewed as a “minimalist” variant of COND that essentially permits an algorithm to compare the relative weights of two items under D (and to draw random samples from D , by setting $S = [N]$).
- **ICOND oracle:** We define an ICOND (short for “interval-cond”) oracle for D as a restricted version of COND_D that only accepts input sets S which are intervals $S = [a, b] = \{a, a+1, \dots, b\}$ for some $a \leq b \in [N]$ (note that taking $a = 1, b = N$ this provides the power of a SAMP_D oracle). This is a natural restriction on COND queries in settings where the N points are endowed with a total order.

1.3 Our results

We give a detailed study of a range of natural distribution testing problems in the COND model and its variants described above, establishing both upper and lower bounds on their query complexity. Our results show that the ability to do conditional sampling provides a significant amount of power to property testers, enabling $\text{polylog}(N)$ -query, or even constant-query, algorithms for problems whose sample complexities in the standard model are $N^{\Omega(1)}$; see Table 1. While we have considered a variety of distribution testing problems in the COND model, our results are certainly not exhaustive, and many directions remain to be explored; we discuss some of these in Section 9.

1.3.1 Testing distributions over unstructured domains

In this initial work on the COND model our main focus has been on the simplest (and, we think, most fundamental) problems in distribution testing, such as testing whether D is the uniform distribution \mathcal{U} ; testing whether $D = D^*$ for an explicitly provided D^* ; testing whether $D_1 = D_2$ given COND_{D_1} and COND_{D_2} oracles; and estimating the variation distance between D and the uniform distribution. In what follows d_{TV} denotes the variation distance.

Testing uniformity. We give a PCOND_D algorithm that tests whether $D = \mathcal{U}$ versus $d_{\text{TV}}(D, \mathcal{U}) \geq \epsilon$ using $\tilde{O}(1/\epsilon^2)$ calls to PCOND_D , independent of N . We show that this PCOND_D algorithm is nearly optimal by proving that any COND_D tester (which may use arbitrary subsets $S \subseteq [N]$ as its query sets) requires $\Omega(1/\epsilon^2)$ queries for this testing problem.

Problem	Our results	Standard model
Is D uniform?	COND $_D$ $\Omega\left(\frac{1}{\epsilon^2}\right)$	$\Theta\left(\frac{\sqrt{N}}{\epsilon^2}\right)$ [GR00, BFR ⁺ 10, Pan08]
	PCOND $_D$ $\tilde{O}\left(\frac{1}{\epsilon^2}\right)$	
	ICOND $_D$ $\tilde{O}\left(\frac{\log^3 N}{\epsilon^3}\right)$ $\Omega\left(\frac{\log N}{\log \log N}\right)$	
Is $D = D^*$ for a known D^* ?	COND $_D$ $\tilde{O}\left(\frac{1}{\epsilon^4}\right)$	$\tilde{\Theta}\left(\frac{\sqrt{N}}{\epsilon^2}\right)$ [BFF ⁺ 01, Pan08]
	PCOND $_D$ $\tilde{O}\left(\frac{\log^4 N}{\epsilon^4}\right)$ $\Omega\left(\sqrt{\frac{\log N}{\log \log N}}\right)$	
Are D_1, D_2 (both unknown) equivalent?	COND $_{D_1, D_2}$ $\tilde{O}\left(\frac{\log^5 N}{\epsilon^4}\right)$	$\tilde{O}\left(\frac{N^{2/3}}{\epsilon^{8/3}}\right)$ [BFR ⁺ 10]
	PCOND $_{D_1, D_2}$ $\tilde{O}\left(\frac{\log^6 N}{\epsilon^{21}}\right)$	$\Omega(N^{2/3})$ [BFR ⁺ 10, Val11]
How far is D from uniform?	PCOND $_D$ $\tilde{O}\left(\frac{1}{\epsilon^{20}}\right)$	$O\left(\frac{1}{\epsilon^2} \frac{N}{\log N}\right)$ [VV11, VV10b] $\Omega\left(\frac{N}{\log N}\right)$ [VV11, VV10a]

Table 1: Comparison between the COND model and the standard model on a variety of distribution testing problems over $[N]$. The upper bounds for the first three problems are for testing whether the property holds (i.e. $d_{TV} = 0$) versus $d_{TV} \geq \epsilon$, and for the last problem the upper bound is for estimating the distance to uniformity to within an additive $\pm\epsilon$.

Testing equivalence to a known distribution. As described above, for the simple problem of testing uniformity we have an essentially optimal PCOND testing algorithm and a matching lower bound. Given these results it is natural to turn to the more challenging question of testing whether D (accessible via a PCOND or COND oracle) is equivalent to D^* , where D^* is an arbitrary “known” distribution over $[N]$ that is explicitly provided to the testing algorithm at no cost (say as a vector $(D^*(1), \dots, D^*(N))$ of probability values). For this “known D^* ” problem, we give a PCOND $_D$ algorithm testing whether $D = D^*$ versus $d_{TV}(D, D^*) \geq \epsilon$ using $\tilde{O}((\log N)^4/\epsilon^4)$ queries. We further show that the $(\log N)^{\Omega(1)}$ query complexity of our PCOND $_D$ algorithm is inherent in the problem, by proving that any PCOND $_D$ algorithm for this problem must use $\sqrt{\log(N)/\log \log(N)}$ queries for constant ϵ .

Given these $(\log N)^{\Theta(1)}$ upper and lower bounds on the query complexity of PCOND $_D$ -testing equivalence to a known distribution, it is natural to ask whether the full COND $_D$ oracle provides more power for this problem. We show that this is indeed the case, by giving a $\tilde{O}(1/\epsilon^4)$ -query algorithm (independent of N) that uses unrestricted COND $_D$ queries.

Testing equivalence between two unknown distributions. We next consider the more challenging problem of testing whether two unknown distributions D_1, D_2 over $[N]$ (available via COND $_{D_1}$ and COND $_{D_2}$ oracles) are identical versus ϵ -far. We give two very different algorithms for this problem. The first uses PCOND oracles and has query complexity $\tilde{O}((\log N)^6/\epsilon^{21})$, while the second uses COND oracles and has query complexity $\tilde{O}((\log N)^5/\epsilon^4)$. We believe that the proof technique of the second algorithm is of independent interest, since it shows how a COND $_D$ oracle

can efficiently simulate an “approximate EVAL_D oracle.” (An EVAL_D oracle takes as input a point $i \in [N]$ and outputs the probability mass $D(i)$ that D puts on i ; we briefly explain our notion of approximating such an oracle in Subsection 1.3.3.)

Estimating the distance to uniformity. We also consider the problem of estimating the variation distance between D and the uniform distribution \mathcal{U} over $[N]$, to within an additive error of $\pm\epsilon$. In the standard SAMP_D model this is known to be a very difficult problem, with an $\Omega(N/\log N)$ lower bound established in [VV11, VV10a]. In contrast, we give a PCOND_D algorithm that makes only $\tilde{O}(1/\epsilon^{20})$ queries, independent of N .

1.3.2 Testing distributions over structured domains

In the final portion of the paper we view the domain $[N]$ as an ordered set $1 \leq \dots \leq N$. (Note that in all the testing problems and results described previously, the domain could just as well have been viewed as an unstructured set of abstract points x_1, \dots, x_N .) With this perspective it is natural to consider ICOND algorithms.

We give an $\tilde{O}((\log N)^3/\epsilon^3)$ -query ICOND_D algorithm for testing whether D is uniform versus ϵ -far from uniform. We show that a $(\log N)^{\Omega(1)}$ query complexity is inherent for uniformity testing using ICOND_D , by proving an $\Omega(\log N/\log \log N)$ -query ICOND_D lower bound.

1.3.3 A high-level discussion of our algorithms

To maintain focus here we describe only the ideas behind our algorithms; intuition for each of our lower bounds can be found in an informal discussion preceding the formal proof, see the beginnings of Sections 4.2, 5.2, and 8.2. As can be seen in the following discussion, our algorithms share some common themes, though each has its own unique idea/technique, which we emphasize below.

Our simplest testing algorithm is the algorithm for **testing whether D is uniform** over $[N]$ (using PCOND_D queries). The algorithm is based on the observation that if a distribution is ϵ -far from uniform, then the total weight (according to D) of points $y \in [N]$ for which $D(y) \geq (1 + \Omega(\epsilon))/N$ is $\Omega(\epsilon)$, and the fraction of points $x \in [N]$ for which $D(x) \leq (1 - \Omega(\epsilon))/N$ is $\Omega(\epsilon)$. If we obtain such a pair of points (x, y) , then we can detect this deviation from uniformity by performing $\Theta(1/\epsilon^2)$ PCOND_D queries on the pair. Such a pair can be obtained with high probability by making $\Theta(1/\epsilon)$ SAMP_D queries (so as to obtain y) as well as selecting $\Theta(1/\epsilon)$ points uniformly (so as to obtain x). This approach yields an algorithm whose complexity grows like $1/\epsilon^4$. To actually get an algorithm with query complexity $\tilde{O}(1/\epsilon^2)$ (which, as our lower bound shows, is tight), a slightly more refined approach is applied.

When we take the next step to **testing equality to an arbitrary (but fully specified) distribution D^*** , the abovementioned observation generalizes so as to imply that if we sample $\Theta(1/\epsilon)$ points from D and $\Theta(1/\epsilon)$ from D^* , then with high probability we shall obtain a pair of points (x, y) such that $D(x)/D(y)$ differs by at least $(1 \pm \Omega(\epsilon))$ from $D^*(x)/D^*(y)$. Unfortunately, this cannot necessarily be detected by a small number of PCOND_D queries since (as opposed to the uniform case), $D^*(x)/D^*(y)$ may be very large or very small. However, we show that by sampling from both D and D^* and allowing the number of samples to grow with $\log N$, with high probability

we either obtain a pair of points as described above for which $D^*(x)/D^*(y)$ is a constant, or we detect that for some set of points B we have that $|D(B) - D^*(B)|$ is relatively large.⁵

As noted previously, we prove a lower bound showing that a polynomial dependence on $\log N$ is unavoidable if only PCOND_D queries (in addition to standard sampling) are allowed. To obtain our more efficient $\text{poly}(1/\epsilon)$ -queries algorithm, which uses more general COND_D queries, we extend the observation from the uniform case in a different way. Specifically, rather than comparing the relative weight of pairs of points, we compare the relative weight of pairs in which one element is a point and the other is a subset of points. Roughly speaking, we show how points can be paired with subsets of points of comparable weight (according to D^*) such that the following holds. If D is far from D^* , then by taking $\tilde{O}(1/\epsilon)$ samples from D and selecting subsets of points in an appropriate manner (depending on D^*), we can obtain (with high probability) a point x and a subset Y such that $D(x)/D(Y)$ differs significantly from $D^*(X)/D^*(y)$ and $D^*(x)/D^*(Y)$ is a constant.

In our next step, to **testing equality between two unknown distributions** D_1 and D_2 , we need to cope with the fact that we no longer “have a hold” on a known distribution. Our PCOND algorithm can be viewed as creating such a hold in the following sense. By sampling from D_1 we obtain (with high probability) a (relatively small) set of points R that *cover* the distribution D_1 . By “covering” we mean that except for a subset having small weight according to D_1 , all points y in $[N]$ have a *representative* $r \in R$, i.e. a point r such that $D_1(y)$ is close to $D_1(r)$. We then show that if D_2 is far from D_1 , then one of the following must hold: (1) There is relatively large weight, either according to D_1 or according to D_2 , on points y such that for some $r \in R$ we have that $D_1(y)$ is close to $D_1(r)$ but $D_2(y)$ is not sufficiently close to $D_2(r)$; (2) There exists a point $r \in R$ such that the set of points y for which $D_1(y)$ is close to $D_1(r)$ has significantly different weight according to D_2 as compared to D_1 . We note that this algorithm can be viewed as a variant of the PCOND algorithm for the case when one of the distributions is known (where the “buckets” B , which were defined by D^* in that algorithm (and were disjoint), are now defined by the points in R (and are not necessarily disjoint)).

As noted previously, our (general) COND algorithm for testing the equality of two (unknown) distributions is based on a subroutine that estimates $D(x)$ (to within $(1 \pm O(\epsilon))$) for a given point x given access to COND_D . Obtaining such an estimate for *every* $x \in [N]$ cannot be done efficiently for some distributions.⁶ However, we show that if we allow the algorithm to output **UNKNOWN** on some subset of points with total weight $O(\epsilon)$, then the relaxed task can be performed using $\text{poly}(\log N, 1/\epsilon)$ queries, by performing a kind of randomized binary search “with exceptions”. This relaxed version, which we refer to as an *approximate EVAL oracle*, suffices for our needs in distinguishing between the case that D_1 and D_2 are the same distribution and the case in which they are far from each other. It is possible that this procedure will be useful for other tasks as well.

The algorithm for **estimating the distance to uniformity** (which uses $\text{poly}(1/\epsilon)$ PCOND_D queries) is based on a subroutine for finding a *reference point* x together with an estimate $\hat{D}(x)$ of $D(x)$. A reference point should be such that $D(x)$ is relatively close to $1/N$ (if such a point

⁵Here we use B for “Bucket”, as we consider a bucketing of the points in $[N]$ based on their weight according to D^* . We note that bucketing has been used extensively in the context of testing properties of distributions, see e.g. [BFR⁺10, BFF⁺01].

⁶As an extreme case consider a distribution D for which $D(1) = 1 - \phi$ and $D(2) = \dots = D(n) = \phi/(n - 1)$ for some very small ϕ (which in particular may depend on n), and for which we are interested in estimating $D(2)$. This requires $\Omega(1/\phi)$ queries.

cannot be found then it is evidence that D is very far from uniform). Given a reference point x (together with $\widehat{D}(x)$) it is possible to estimate the distance to uniformity by obtaining (using PCOND queries) estimates of the ratio between $D(x)$ and $D(y)$ for $\text{poly}(1/\epsilon)$ uniformly selected points y . The procedure for finding a reference point x together with $\widehat{D}(x)$ is based on estimating both the weight and the size of a subset of points y such that $D(y)$ is close to $D(x)$. The procedure shares a common subroutine, ESTIMATE-NEIGHBORHOOD, with the PCOND algorithm for testing equivalence between two unknown distributions.

Finally, the ICOND $_D$ algorithm for testing uniformity is based on a version of the approximate EVAL oracle mentioned previously, which on one hand uses only ICOND $_D$ (rather than general COND $_D$) queries, and on the other hand exploits the fact that we are dealing with the uniform distribution rather than an arbitrary distribution.

1.4 The independent work [CFG13]

In what follows we discuss the work of Chakraborty et al. [CFG13], which was done independently from our work, and was recently accepted to the ITC conference (so that we learned about its existence only a few days before the STOC submission deadline). Chakraborty et al. [CFG13] propose essentially the same COND model that we propose, differing only in what happens on query sets S such that $D(S) = 0$. In our model such a query causes the COND oracle and algorithm to return FAIL, while in their model such a query returns a uniform random $i \in S$. They present the following results.

- An (adaptive) algorithm for testing uniformity that performs $\text{poly}(1/\epsilon)$ queries.⁷ The sets on which the algorithm performs COND queries are of size linear in $1/\epsilon$. Recall that our algorithm for this problem performs $\tilde{O}(1/\epsilon^2)$ PCOND queries and that we show that every algorithm must perform $\Omega(1/\epsilon^2)$ queries (when there is no restriction on the types of queries). We note that their analysis uses the same observation that ours does regarding distributions that are far from uniform (see the discussion in Subsection 1.3.3), but exploits it in a different manner.

They also give a non-adaptive algorithm for this problem that performs $\text{poly}(\log N, 1/\epsilon)$ COND queries and show that $\Omega(\log \log N)$ is a lower bound on the necessary number of queries for non-adaptive algorithms.

- An (adaptive) algorithm for testing whether D is equivalent to a specified distribution D^* using $\text{poly}(\log^* N, 1/\epsilon)$ COND queries. Recall that we give an algorithm for this problem that performs $\tilde{O}(1/\epsilon^4)$ COND queries.

They also give a non-adaptive algorithm for this problem that performs $\text{poly}(\log N, 1/\epsilon)$ COND queries.

- An (adaptive) algorithm for testing any label-invariant (i.e., invariant under permutations of the domain) property that performs $\text{poly}(\log N, 1/\epsilon)$ COND queries. As noted in [CFG13], this in particular implies an algorithm with this complexity for estimating the distance to

⁷The precise polynomial is not specified – we believe it is roughly $1/\epsilon^4$ as it follows from an application of the identity tester of [BFF⁺01] with distance $\Theta(\epsilon^2)$ on a domain of size $O(1/\epsilon)$.

uniformity. Recall that we give an algorithm for this estimation problem that performs $\text{poly}(1/\epsilon)$ PCOND queries.

The algorithm for testing any label-invariant property is based on learning a certain approximation of the distribution D and in this process defining some sort of approximate EVAL oracle. To the best of our understanding, our notion of an approximate EVAL oracle (which is used to obtain one or our results for testing equivalence between two unknown distributions) is quite different.

They also show that there exists a label-invariant property for which any adaptive algorithm must perform $\Omega(\sqrt{\log \log N})$ COND queries.

- Finally they show that there exist general properties that require $\Omega(n)$ COND queries.

2 Preliminaries

2.1 Definitions

Throughout the paper we shall work with discrete distributions over an N -element set whose elements are denoted $\{1, \dots, N\}$; we write $[N]$ to denote $\{1, \dots, N\}$ and $[a, b]$ to denote $\{a, \dots, b\}$. For a distribution D over $[N]$ we write $D(i)$ to denote the probability of i under D , and for $S \subseteq [N]$ we write $D(S)$ to denote $\sum_{i \in S} D(i)$. For $S \subseteq [N]$ such that $D(S) > 0$ we write D_S to denote the conditional distribution of D restricted to S , so $D_S(i) = \frac{D(i)}{D(S)}$ for $i \in S$ and $D_S(i) = 0$ for $i \notin S$.

As is standard in property testing of distributions, throughout this work we measure the distance between two distributions D_1 and D_2 using the *total variation distance*:

$$d_{\text{TV}}(D_1, D_2) \stackrel{\text{def}}{=} \frac{1}{2} \|D_1 - D_2\|_1 = \frac{1}{2} \sum_{i \in [N]} |D_1(i) - D_2(i)| = \max_{S \subseteq [N]} |D_1(S) - D_2(S)|.$$

We may view a *property* \mathcal{P} of distributions over $[N]$ as a subset of all distributions over $[N]$ (consisting of all distributions that have the property). The distance from D to a property \mathcal{P} , denoted $d_{\text{TV}}(D, \mathcal{P})$, is defined as $\inf_{D' \in \mathcal{P}} \{d_{\text{TV}}(D, D')\}$.

We define testing algorithms for properties of distributions over $[N]$ as follows:

Definition 2 Let \mathcal{P} be a property of distributions over $[N]$. Let ORACLE_D be some type of oracle which provides access to D . A $q(\epsilon, N)$ -query ORACLE testing algorithm for \mathcal{P} is an algorithm T which is given ϵ, N as input parameters and oracle access to an ORACLE_D oracle. For any distribution D over $[N]$ algorithm T makes at most $q(\epsilon, N)$ calls to ORACLE_D , and:

- if $D \in \mathcal{P}$ then with probability at least $2/3$ algorithm T outputs ACCEPT;
- if $d_{\text{TV}}(D, \mathcal{P}) \geq \epsilon$ then with probability at least $2/3$ algorithm T outputs REJECT.

This definition can easily be extended to cover situations in which there are two “unknown” distributions D_1, D_2 that are accessible via ORACLE_{D_1} and ORACLE_{D_2} oracles. In particular we shall consider algorithms for testing whether $D_1 = D_2$ versus $d_{\text{TV}}(D_1, D_2)$ in such a setting. We sometimes write T^{ORACLE_D} to indicate that T has access to ORACLE_D .

2.2 Useful tools

On several occasions we will use the *data processing inequality for variation distance*. This fundamental result says that for any two distributions D, D' , applying any (possibly randomized) function to D and D' can never increase their statistical distance; see e.g. part (iv) of Lemma 2 of [Rey11] for a proof of this lemma.

Lemma 1 (Data Processing Inequality for Total Variation Distance) *Let D, D' be two distributions over a domain Ω . Fix any randomized function⁸ F on Ω , and let $F(D)$ be the distribution such that a draw from $F(D)$ is obtained by drawing independently x from D and f from F and then outputting $f(x)$ (likewise for $F(D')$). Then we have*

$$d_{\text{TV}}(F(D), F(D')) \leq d_{\text{TV}}(D, D').$$

We next give several variants of Chernoff bounds (see e.g. Chapter 4 of [MR95]).

Theorem 1 *Let Y_1, \dots, Y_m be m independent random variables that take on values in $[0, 1]$, where $E[Y_i] = p_i$, and $\sum_{i=1}^m p_i = P$. For any $\gamma \in (0, 1]$ we have*

$$\text{(additive bound)} \quad \Pr \left[\sum_{i=1}^m Y_i > P + \gamma m \right], \Pr \left[\sum_{i=1}^m Y_i < P - \gamma m \right] \leq \exp(-2\gamma^2 m) \quad (1)$$

$$\text{(multiplicative bound)} \quad \Pr \left[\sum_{i=1}^m Y_i > (1 + \gamma)P \right] < \exp(-\gamma^2 P/3) \quad (2)$$

and

$$\text{(multiplicative bound)} \quad \Pr \left[\sum_{i=1}^m Y_i < (1 - \gamma)P \right] < \exp(-\gamma^2 P/2). \quad (3)$$

The bound in Equation (2) is derived from the following more general bound, which holds from any $\gamma > 0$:

$$\Pr \left[\sum_{i=1}^m Y_i > (1 + \gamma)P \right] \leq \left(\frac{e^\gamma}{(1 + \gamma)^{1+\gamma}} \right)^P, \quad (4)$$

and which also implies that for any $B > 2eP$,

$$\Pr \left[\sum_{i=1}^m Y_i > B \right] \leq 2^{-B}. \quad (5)$$

⁸Which can be seen as a distribution over functions over Ω .

The following extension of the multiplicative bound is useful when we only have upper and/or lower bounds on P (see Exercise 1.1 of [DP09]):

Corollary 2 *In the setting of Theorem 1 suppose that $P_L \leq P \leq P_H$. Then for any $\gamma \in (0, 1]$, we have*

$$\Pr \left[\sum_{i=1}^m Y_i > (1 + \gamma)P_H \right] < \exp(-\gamma^2 P_H/3) \quad (6)$$

$$\Pr \left[\sum_{i=1}^m Y_i < (1 - \gamma)P_L \right] < \exp(-\gamma^2 P_L/2) \quad (7)$$

We will also use the following corollary of Theorem 1:

Corollary 3 *Let $0 \leq w_1, \dots, w_m \in \mathbb{R}$ be such that $w_i \leq \kappa$ for all $i \in [m]$ where $\kappa \in (0, 1]$. Let X_1, \dots, X_m be i.i.d. Bernoulli random variables with $\Pr[X_i = 1] = 1/2$ for all i , and let $X = \sum_{i=1}^m w_i X_i$ and $W = \sum_{i=1}^m w_i$. For any $\gamma \in (0, 1]$,*

$$\Pr \left[X > (1 + \gamma) \frac{W}{2} \right] < \exp \left(-\gamma^2 \frac{W}{6\kappa} \right) \quad \text{and} \quad \Pr \left[X < (1 - \gamma) \frac{W}{2} \right] < \exp \left(-\gamma^2 \frac{W}{4\kappa} \right),$$

and for any $B > e \cdot W$,

$$\Pr[X > B] < 2^{-B/\kappa}.$$

Proof: Let $w'_i = w_i/\kappa$ (so that $w'_i \in [0, 1]$), let $W' = \sum_{i=1}^m w'_i = W/\kappa$, and for each $i \in [m]$ let $Y_i = w'_i X_i$, so that Y_i takes on values in $[0, 1]$ and $E[Y_i] = w'_i/2$. Let $X' = \sum_{i=1}^m w'_i X_i = \sum_{i=1}^m Y_i$, so that $E[X'] = W'/2$. By the definitions of W' and X' and by Equation (2), for any $\gamma \in (0, 1]$,

$$\Pr \left[X > (1 + \gamma) \frac{W}{2} \right] = \Pr \left[X' > (1 + \gamma) \frac{W'}{2} \right] < \exp \left(-\gamma^2 \frac{W'}{6} \right) = \exp \left(-\gamma^2 \frac{W}{6\kappa} \right), \quad (8)$$

and similarly by Equation (3)

$$\Pr \left[X < (1 - \gamma) \frac{W}{2} \right] < \exp \left(-\gamma^2 \frac{W}{4\kappa} \right). \quad (9)$$

For $B > e \cdot W = 2e \cdot W/2$ we apply Equation (5) and get

$$\Pr[X > B] = \Pr[X' > B/\kappa] < 2^{-B/\kappa}, \quad (10)$$

as claimed. ■

3 Some useful procedures

In this section we describe some procedures that will be used by our algorithms.

3.1 The procedure COMPARE

We start by describing a procedure that estimates the ratio between the weights of two disjoint sets of points by performing COND queries on the union of the sets. In the special case when each set is of size one, the queries performed are PCOND queries.

Algorithm 1: COMPARE

Input: COND query access to a distribution D over $[N]$, disjoint subsets $X, Y \subset [N]$, parameters $\eta \in (0, 1]$, $K \geq 1$, and $\delta \in (0, 1/2]$.

1. Perform $\Theta\left(\frac{K \log(1/\delta)}{\eta^2}\right)$ COND_D queries on the set $S = X \cup Y$, and let $\hat{\mu}$ be the fraction of times that a point $y \in Y$ is returned.
 2. If $\hat{\mu} < \frac{2}{3} \cdot \frac{1}{K+1}$, then return **Low**.
 3. Else, if $1 - \hat{\mu} < \frac{2}{3} \cdot \frac{1}{K+1}$, then return **High**.
 4. Else return $\rho = \frac{\hat{\mu}}{1-\hat{\mu}}$.
-

Lemma 2 *Given as input two disjoint subsets of points X, Y together with parameters $\eta \in (0, 1]$, $K \geq 1$, and $\delta \in (0, 1/2]$, as well as COND query access to a distribution D , the procedure COMPARE (Algorithm 1) either outputs a value $\rho > 0$ or outputs **High** or **Low**, and satisfies the following:*

1. *If $D(X)/K \leq D(Y) \leq K \cdot D(X)$ then with probability at least $1 - \delta$ the procedure outputs a value $\rho \in [1 - \eta, 1 + \eta]D(Y)/D(X)$;*
2. *If $D(Y) > K \cdot D(X)$ then with probability at least $1 - \delta$ the procedure outputs either **High** or a value $\rho \in [1 - \eta, 1 + \eta]D(Y)/D(X)$;*
3. *If $D(Y) < D(X)/K$ then with probability at least $1 - \delta$ the procedure outputs either **Low** or a value $\rho \in [1 - \eta, 1 + \eta]D(Y)/D(X)$.*

The procedure performs $O\left(\frac{K \log(1/\delta)}{\eta^2}\right)$ COND queries on the set $X \cup Y$.

Proof: The bound on the number of queries performed by the algorithm follows directly from the description of the algorithm, and hence we turn to establish its correctness.

Let $w(X) = \frac{D(X)}{D(X)+D(Y)}$ and let $w(Y) = \frac{D(Y)}{D(X)+D(Y)}$. Observe that $\frac{w(Y)}{w(X)} = \frac{D(Y)}{D(X)}$ and that for $\hat{\mu}$ as defined in Line 1 of the algorithm, $E[\hat{\mu}] = w(Y)$ and $E[1 - \hat{\mu}] = w(X)$. Also observe that for any $B \geq 1$, if $D(Y) \geq D(X)/B$, then $w(Y) \geq \frac{1}{B+1}$ and if $D(Y) \leq B \cdot D(X)$, then $w(X) \geq \frac{1}{B+1}$.

Let E_1 be the event that $\hat{\mu} \in [1 - \eta/3, 1 + \eta/3]w(Y)$ and let E_2 be the event that $(1 - \hat{\mu}) \in [1 - \eta/3, 1 + \eta/3]w(X)$. Given the number of COND queries performed on the set $X \cup Y$, by applying a multiplicative Chernoff bound (see Theorem 1), if $w(Y) \geq \frac{1}{4K}$ then with probability at least $1 - \delta/2$ the event E_1 holds, and if $w(X) \geq \frac{1}{4K}$, then with probability at least $1 - \delta/2$ the event E_2 holds. We next consider the three cases in the lemma statement.

1. If $D(X)/K \leq D(Y) \leq KD(X)$, then by the discussion above, $w(Y) \geq \frac{1}{K+1}$, $w(X) \geq \frac{1}{K+1}$, and with probability at least $1 - \delta$ we have that $\hat{\mu} \in [1 - \eta/3, 1 + \eta/3]w(Y)$ and $(1 - \hat{\mu}) \in [1 - \eta/3, 1 + \eta/3]w(X)$. Conditioned on these bounds holding,

$$\hat{\mu} \geq \frac{1 - \eta/3}{K + 1} \geq \frac{2}{3} \cdot \frac{1}{K + 1} \quad \text{and} \quad 1 - \hat{\mu} \geq \frac{2}{3} \cdot \frac{1}{K + 1}.$$

It follows that the procedure outputs a value $\rho = \frac{\hat{\mu}}{1 - \hat{\mu}} \in [1 - \eta, 1 + \eta] \frac{w(Y)}{w(X)}$ as required by Item 1.

2. If $D(Y) > K \cdot D(X)$, then we consider two subcases.

- (a) If $D(Y) > 3K \cdot D(X)$, then $w(X) < \frac{1}{3K+1}$, so that by a multiplicative Chernoff bound (stated in Corollary 2), with probability at least $1 - \delta$ we have that

$$1 - \hat{\mu} < \frac{1 + \eta/3}{3K + 1} \leq \frac{4}{3} \cdot \frac{1}{3K + 1} \leq \frac{2}{3} \cdot \frac{1}{K + 1},$$

causing the algorithm to output **High**. Thus Item 2 is established for this subcase.

- (b) If $K \cdot D(X) < D(Y) \leq 3K \cdot D(X)$, then $w(X) \geq \frac{1}{3K+1}$ and $w(Y) \geq \frac{1}{2}$, so that the events E_1 and E_2 both hold with probability at least $1 - \delta$. Assume that these events in fact hold. This implies that $\hat{\mu} \geq \frac{1 - \eta/3}{2} \geq \frac{2}{3} \cdot \frac{1}{K+1}$, and the algorithm either outputs **High** or outputs $\rho = \frac{\hat{\mu}}{1 - \hat{\mu}} \in [1 - \eta, 1 + \eta] \frac{w(Y)}{w(X)}$, so Item 2 is established for this subcase as well.

3. If $D(Y) < D(X)/K$, so that $D(X) > K \cdot D(Y)$, then the exact same arguments are applied as in the previous case, just switching the roles of Y and X and the roles of $\hat{\mu}$ and $1 - \hat{\mu}$ so as to establish Item 3.

We have thus established all items in the lemma. ■

3.2 The procedure ESTIMATE-NEIGHBORHOOD

In this subsection we describe a procedure that, given a point x , provides an estimate of the weight of a set of points y such that $D(y)$ is similar to $D(x)$. In order to specify the behavior of the procedure more precisely, we introduce the following notation. For a distribution D over $[N]$, a point $x \in [N]$ and a parameter $\gamma \in [0, 1]$, let

$$U_\gamma^D(x) \stackrel{\text{def}}{=} \left\{ y \in [N] : \frac{1}{1 + \gamma} D(x) \leq D(y) \leq (1 + \gamma) D(x) \right\} \quad (11)$$

denote the set of points whose weight is “ γ -close” to the weight of x . If we take a sample of points distributed according to D , then the expected fraction of these points that belong to $U_\gamma^D(x)$ is $D(U_\gamma^D(x))$. If this value is not too small, then the actual fraction in the sample is close to the expected value. Hence, if we could efficiently determine for any given point y whether or not it belongs to $U_\gamma^D(x)$, then we could obtain a good estimate of $D(U_\gamma^D(x))$. The difficulty is that it is not possible to perform this task efficiently for “boundary” points y such that $D(y)$ is very close to $(1 + \gamma)D(x)$ or to $\frac{1}{1 + \gamma}D(x)$. However, for our purposes, it is not important that we obtain the weight and size of $U_\gamma^D(x)$ for a specific γ , but rather it suffices to do so for γ in a given range, as stated in the next lemma.

Lemma 3 *Given as input a point x together with parameters $\kappa, \beta, \eta, \delta \in (0, 1/2]$ as well as PCOND query access to a distribution D , the procedure ESTIMATE-NEIGHBORHOOD (Algorithm 2) outputs a pair $(\hat{w}, \alpha) \in [0, 1] \times (\kappa, 2\kappa)$ such that α is uniformly distributed in $\{\kappa + i\theta\}_{i=0}^{\kappa/\theta-1}$ for $\theta = \frac{\kappa\eta\beta\delta}{64}$, and such that the following holds:*

1. *If $D(U_\alpha^D(x)) \geq \beta$, then with probability at least $1 - \delta$ we have $\hat{w} \in [1 - \eta, 1 + \eta] \cdot D(U_\alpha^D(x))$, and $D(U_{\alpha+\theta}^D(x) \setminus U_\alpha^D(x)) \leq \eta\beta/16$;*
2. *If $D(U_\alpha^D(x)) < \beta$, then with probability at least $1 - \delta$ we have $\hat{w} \leq (1 + \eta) \cdot \beta$, and $D(U_{\alpha+\theta}^D(x) \setminus U_\alpha^D(x)) \leq \eta\beta/16$.*

The number of PCOND queries performed by the procedure is $O\left(\frac{\log(1/\delta) \cdot \log(\log(1/\delta)/(\beta\eta^2))}{\kappa^2 \eta^4 \beta^3 \delta^2}\right)$.

Algorithm 2: ESTIMATE-NEIGHBORHOOD

Input: PCOND query access to a distribution D over $[N]$, a point $x \in [N]$ and parameters

- $\kappa, \beta, \eta, \delta \in (0, 1/2]$
- 1: Set $\theta = \frac{\kappa\eta\beta\delta}{64}$ and $r = \frac{\kappa}{\theta} = \frac{64}{\eta\beta\delta}$.
 - 2: Select a value $\alpha \in \{\kappa + i\theta\}_{i=0}^{r-1}$ uniformly at random.
 - 3: Call the SAMP_D oracle $\Theta(\log(1/\delta)/(\beta\eta^2))$ times and let S be the set of points obtained.
 - 4: For each point y in S call $\text{COMPARE}_D(\{x\}, \{y\}, \theta/4, 4, \delta/(4|S|))$ (if a point y appears more than once in S , then COMPARE is called only once on y).
 - 5: Let \hat{w} be the fraction of occurrences of points y in S for which COMPARE returned a value $\rho(y) \in [1/(1 + \alpha + \theta/2), (1 + \alpha + \theta/2)]$. (That is, S is viewed as a multiset.)
 - 6: Return (\hat{w}, α) .
-

Proof of Lemma 3: The number of PCOND queries performed by ESTIMATE-NEIGHBORHOOD is the size of S times the number of PCOND queries performed in each call to COMPARE. By the setting of the parameters in the calls to COMPARE, the total number of PCOND queries is $O\left(\frac{|S| \cdot \log |S| / \delta}{\theta^2}\right) = O\left(\frac{\log(1/\delta) \cdot \log(\log(1/\delta)/(\beta\eta^2))}{\kappa^2 \eta^4 \beta^3 \delta^2}\right)$. We now turn to establishing the correctness of the procedure.

Since D and x are fixed, in what follows we shall use the shorthand U_γ for $U_\gamma^D(x)$. For $\alpha \in \{\kappa + i\theta\}_{i=0}^{r-1}$, let $\Delta_\alpha \stackrel{\text{def}}{=} U_{\alpha+\theta} \setminus U_\alpha$. We next define several “desirable” events. In all that follows we view S as a multiset.

1. Let E_1 be the event that $D(\Delta_\alpha) \leq 4/(\delta r)$. Since there are r disjoint sets Δ_α for $\alpha \in \{\kappa + i\theta\}_{i=0}^{r-1}$, the probability that E_1 occurs (taken over the uniform choice of α) is at least $1 - \delta/4$. From this point on we fix α and assume E_1 holds.
2. The event E_2 is that $|S \cap \Delta_\alpha|/|S| \leq 8/(\delta r)$ (that is, at most twice the upper bound on the expected value). By applying the multiplicative Chernoff bound using the fact that $|S| = \Theta(\log(1/\delta)/(\beta\eta^2)) = \Omega(\log(1/\delta) \cdot (\delta r))$, we have that $\Pr_S[E_2] \geq 1 - \delta/4$.

3. The event E_3 is defined as follows: If $D(U_\alpha) \geq \beta$, then $|S \cap U_\alpha|/|S| \in [1 - \eta/2, 1 + \eta/2] \cdot D(U_\alpha)$, and if $D(U_\alpha) < \beta$, then $|S \cap U_\alpha|/|S| < (1 + \eta/2) \cdot \beta$. Once again applying the multiplicative Chernoff bound (for both cases) and using that fact that $|S| = \Theta(\log(1/\delta)/(\beta\eta^2))$, we have that $\Pr_S[E_3] \geq 1 - \delta/4$.
4. Let E_4 be the event that all calls to COMPARE return an output as specified in Lemma 2. Given the setting of the confidence parameter in the calls to COMPARE we have that $\Pr[E_4] \geq 1 - \delta/4$ as well.

Assume from this point on that events E_1 through E_4 all hold where this occurs with probability at least $1 - \delta$. By the definition of Δ_α and E_1 we have that $D(U_{\alpha+\theta} \setminus U_\alpha) \leq 4/(\delta r) = \eta\beta/16$, as required (in both items of the lemma). Let T be the (multi-)subset of points y in S for which COMPARE returned a value $\rho(y) \in [1/(1 + \alpha + \theta/2), (1 + \alpha + \theta/2)]$ (so that \hat{w} , as defined in the algorithm, equals $|T|/|S|$). Note first that conditioned on E_4 we have that for every $y \in U_{2\kappa}$ it holds that the output of COMPARE when called on $\{x\}$ and $\{y\}$, denoted $\rho(y)$, satisfies $\rho(y) \in [1 - \theta/4, 1 + \theta/4](D(y)/D(x))$, while for $y \notin U_{2\kappa}$ either COMPARE outputs **High** or **Low** or it outputs a value $\rho(y) \in [1 - \theta/4, 1 + \theta/4](D(y)/D(x))$. This implies that if $y \in U_\alpha$, then $\rho(y) \leq (1 + \alpha) \cdot (1 + \theta/4) \leq 1 + \alpha + \theta/2$ and $\rho(y) \geq (1 + \alpha)^{-1} \cdot (1 - \theta/4) \geq (1 + \alpha + \theta/2)^{-1}$, so that $S \cap U_\alpha \subseteq T$. On the other hand, if $y \notin U_{\alpha+\theta}$ then either $\rho(y) > (1 + \alpha + \theta) \cdot (1 - \theta/4) \geq 1 + \alpha + \theta/2$ or $\rho(y) < (1 + \alpha + \theta)^{-1} \cdot (1 + \theta/4) \leq (1 + \alpha + \theta/2)^{-1}$ so that $T \subseteq S \cap U_{\alpha+\theta}$. Combining the two we have:

$$S \cap U_\alpha \subseteq T \subseteq S \cap U_{\alpha+\theta}. \quad (12)$$

Recalling that $\hat{w} = \frac{|T|}{|S|}$, the left-hand side of Equation (12) implies that

$$\hat{w} \geq \frac{|S \cap U_\alpha|}{|S|}, \quad (13)$$

and by E_1 and E_2 , the right-hand-side of Equation (12) implies that

$$\hat{w} \leq \frac{|S \cap U_\alpha|}{|S|} + \frac{8}{\delta r} \leq \frac{|S \cap U_\alpha|}{|S|} + \frac{\beta\eta}{8}. \quad (14)$$

We consider the two cases stated in the lemma:

1. If $D(U_\alpha) \geq \beta$, then by Equation (13), Equation (14) and (the first part of) E_3 , we have that $\hat{w} \in [1 - \eta, 1 + \eta] \cdot D(U_\alpha)$.
2. If $D(U_\alpha) < \beta$, then by Equation (14) and (the second part of) E_3 , we have that $\hat{w} \leq (1 + \eta)\beta$.

The lemma is thus established. \blacksquare

4 Algorithms and lower bounds for testing uniformity

4.1 A $\tilde{O}(1/\epsilon^2)$ -query PCOND algorithm for testing uniformity

In this subsection we present an algorithm $\text{PCOND}_D\text{-TEST-UNIFORM}$ and prove the following theorem:

Theorem 4 $\text{PCOND}_D\text{-TEST-UNIFORM}$ is a $\tilde{O}(1/\epsilon^2)$ -query PCOND_D testing algorithm for uniformity, i.e. it outputs ACCEPT with probability at least $2/3$ if $D = \mathcal{U}$ and outputs REJECT with probability at least $2/3$ if $d_{\text{TV}}(D, \mathcal{U}) \geq \epsilon$.

Intuition. For the sake of intuition we first describe a simpler approach that yields a $\tilde{O}(1/\epsilon^4)$ -query algorithm, and then build on those ideas to obtain our real algorithm with its improved $\tilde{O}(1/\epsilon^2)$ bound. Fix D to be a distribution over $[N]$ that is ϵ -far from uniform. Let

$$H = \left\{ h \in [N] \mid D(h) \geq \frac{1}{N} \right\} \text{ and } L = \left\{ \ell \in [N] \mid D(\ell) < \frac{1}{N} \right\}.$$

It is easy to see that since D is ϵ -far from uniform, we have

$$\sum_{h \in H} \left(D(h) - \frac{1}{N} \right) = \sum_{\ell \in L} \left(\frac{1}{N} - D(\ell) \right) \geq \frac{\epsilon}{2}. \quad (15)$$

From this it is not hard to show that

- (i) many elements of $[N]$ must be “significantly light” in the following sense: Define $L' \subseteq L$ to be $L' = \{ \ell \in [N] \mid D(\ell) < \frac{1}{N} - \frac{\epsilon}{4N} \}$. Then it must be the case that $|L'| \geq (\epsilon/4)N$.
- (ii) D places significant weight on elements that are “significantly heavy” in the following sense: Define $H' \subseteq H$ to be $H' = \{ h \in [N] \mid D(h) \geq \frac{1}{N} + \frac{\epsilon}{4N} \}$. Then it must be the case that $D(H') \geq (\epsilon/4)$.

Using (i) and (ii) it is fairly straightforward to give a $O(1/\epsilon^4)$ -query PCOND_D testing algorithm as follows: we can get a point in L' with high probability by randomly sampling $O(1/\epsilon)$ points uniformly at random from $[N]$, and we can get a point in H' with high probability by drawing $O(1/\epsilon)$ points from SAMP_D . Then at least one of the $O(1/\epsilon^2)$ pairs that have one point from the first sample and one point from the second will have a multiplicative factor difference of $1 + \Omega(\epsilon)$ between the weight under D of the two points, and this can be detected by calling the procedure COMPARE (see Subsection 3.1). Since there are $O(1/\epsilon^2)$ pairs and for each one the invocation of COMPARE uses $\tilde{O}(1/\epsilon^2)$ queries, the overall sample complexity of this simple approach is $\tilde{O}(1/\epsilon^4)$.

Our actual algorithm $\text{PCOND}_D\text{-TEST-UNIFORM}$ for testing uniformity extends the above ideas to get a $\tilde{O}(1/\epsilon^2)$ -query algorithm. More precisely, the algorithm works as follows: it first draws a “reference sample” of $O(1)$ points uniformly from $[N]$. Next, repeatedly for $O(\log \frac{1}{\epsilon})$ iterations, the algorithm draws two other samples, one uniformly from $[N]$ and the other from SAMP_D . (These samples have different sizes at different iterations; intuitively, each iteration is meant to deal with a different “scale” of probability mass that points could have under D .) At each iteration it then uses COMPARE to do comparisons between pairs of elements, one from the reference sample and the other from one of the two other samples. If D is ϵ -far from uniform, then with high probability at some iteration the algorithm will either draw a point from SAMP_D that has “very big” mass under D , or draw a point from the uniform distribution over $[N]$ that has “very small” mass under D , and this will be detected by the comparisons to the reference points. Choosing the sample sizes and parameters for the COMPARE calls carefully at each iteration yields the improved query bound.

Algorithm 3: PCOND_D-TEST-UNIFORM

Input: error parameter $\epsilon > 0$; query access to PCOND_D oracle

- 1: Set $t = \log(\frac{4}{\epsilon}) + 1$.
 - 2: Select $q = \Theta(1)$ points i_1, \dots, i_q independently and uniformly from $[N]$.
 - 3: **for** $j = 1$ to t **do**
 - 4: Call the SAMP_D oracle $s_j = \Theta(2^j \cdot t)$ times to obtain points h_1, \dots, h_{s_j} distributed according to D .
 - 5: Select s_j points $\ell_1, \dots, \ell_{s_j}$ independently and uniformly from $[N]$.
 - 6: **for all** pairs $(x, y) = (i_r, h_{r'})$ and $(x, y) = (i_r, \ell_{r'})$ (where $1 \leq r \leq q, 1 \leq r' \leq s_j$) **do**
 - 7: Call COMPARE_D({ x }, { y }, $\Theta(\epsilon 2^j), 2, \exp^{-\Theta(t)}$).
 - 8: **if** the COMPARE call does not return a value in $[\frac{1}{2} - 2^{j-5} \frac{\epsilon}{4}, \frac{1}{2} + 2^{j-5} \frac{\epsilon}{4}]$ **then**
 - 9: output REJECT (and exit).
 - 10: **end if**
 - 11: **end for**
 - 12: **end for**
 - 13: Output ACCEPT
-

Let m_j denote the number of PCOND_D queries used to run COMPARE_D in a given execution of Line 7 during the j -th iteration of the outer loop. By the setting of the parameters in each such call and Lemma 2, $m_j = O(\frac{t}{\epsilon^2 2^{2j}})$. It is easy to see that the algorithm only performs PCOND_D queries and that the total number of queries that the algorithm performs is

$$O\left(\sum_{j=1}^t q \cdot s_j \cdot m_j\right) = O\left(\sum_{j=1}^t 2^j \log\left(\frac{1}{\epsilon}\right) \cdot \frac{\log(\frac{1}{\epsilon})}{\epsilon^2 2^{2j}}\right) = O\left(\frac{(\log(\frac{1}{\epsilon}))^2}{\epsilon^2}\right).$$

We prove Theorem 4 by arguing completeness and soundness below.

Completeness: Suppose that D is the uniform distribution. Then for any fixed pair of points (x, y) , Lemma 2 implies that the call to COMPARE on { x }, { y } in Line 7 causes the algorithm to output REJECT in Line 9 with probability at most $e^{-\Theta(t)} = \text{poly}(\epsilon)$. By taking a union bound over all $\text{poly}(1/\epsilon)$ pairs of points considered by the algorithm, the algorithm will accept with probability at least $2/3$, as required.

Soundness: Now suppose that D is ϵ -far from uniform (we assume throughout the analysis that $\epsilon = 1/2^k$ for some integer k , which is clearly without loss of generality). We define H, L as above and further partition H and L into “buckets” as follows: for $j = 1, \dots, t - 1 = \log(\frac{4}{\epsilon})$, let

$$H_j \stackrel{\text{def}}{=} \left\{ h \mid \left(1 + 2^{j-1} \cdot \frac{\epsilon}{4}\right) \cdot \frac{1}{N} \leq D(h) < \left(1 + 2^j \cdot \frac{\epsilon}{4}\right) \cdot \frac{1}{N} \right\},$$

and

$$L_j \stackrel{\text{def}}{=} \left\{ \ell \mid \left(1 - 2^j \cdot \frac{\epsilon}{4}\right) \cdot \frac{1}{N} < D(\ell) \leq \left(1 - 2^{j-1} \cdot \frac{\epsilon}{4}\right) \cdot \frac{1}{N} \right\}.$$

Also define

$$H_0 \stackrel{\text{def}}{=} \left\{ h \mid \frac{1}{N} \leq D(h) < \left(1 + \frac{\epsilon}{4}\right) \cdot \frac{1}{N} \right\}, \quad L_0 \stackrel{\text{def}}{=} \left\{ \ell \mid \left(1 - \frac{\epsilon}{4}\right) \cdot \frac{1}{N} < D(\ell) < \frac{1}{N} \right\},$$

and

$$H_t \stackrel{\text{def}}{=} \left\{ h \mid D(h) \geq \frac{2}{N} \right\}, \quad L_t \stackrel{\text{def}}{=} \left\{ \ell \mid D(\ell) \leq \frac{1}{2N} \right\}.$$

First observe that by the definition of H_0 and L_0 , we have

$$\sum_{h \in H_0} \left(D(h) - \frac{1}{N} \right) \leq \frac{\epsilon}{4} \quad \text{and} \quad \sum_{\ell \in L_0} \left(\frac{1}{N} - D(\ell) \right) \leq \frac{\epsilon}{4}.$$

Therefore (by Equation (15)) we have

$$\sum_{j=1}^t \sum_{h \in H_j} \left(D(h) - \frac{1}{N} \right) \geq \frac{\epsilon}{4} \quad \text{and} \quad \sum_{j=1}^t \sum_{\ell \in L_j} \left(\frac{1}{N} - D(\ell) \right) \geq \frac{\epsilon}{4}.$$

This implies that for some $1 \leq j(H) \leq t$, and some $1 \leq j(L) \leq t$ we have

$$\sum_{h \in H_{j(H)}} \left(D(h) - \frac{1}{N} \right) \geq \frac{\epsilon}{4t} \quad \text{and} \quad \sum_{\ell \in L_{j(L)}} \left(\frac{1}{N} - D(\ell) \right) \geq \frac{\epsilon}{4t}. \quad (16)$$

The rest of the analysis is divided into two cases depending on whether $|L| \geq \frac{N}{2}$ or $|H| > \frac{N}{2}$.

Case 1: $|L| \geq \frac{N}{2}$. In this case, with probability at least 99/100, in Line 2 the algorithm will select at least one point $i_r \in L$. We consider two subcases: $j(H) = t$, and $j(H) \leq t - 1$.

- $j(H) = t$: In this subcase, by Equation (16) we have that $\sum_{h \in H_{j(H)}} D(h) \geq \frac{\epsilon}{4t}$. This implies that when $j = j(H) = t = \log(\frac{4}{\epsilon}) + 1$, so that $s_j = s_t = \Theta(\frac{t}{\epsilon})$, with probability at least 99/100 the algorithm selects a point $h_{r'} \in H_t$ in Line 4. Assume that indeed such a point $h_{r'}$ is selected. Since $D(h_{r'}) \geq \frac{2}{N}$, while $D(i_r) < \frac{1}{N}$, Lemma 2 implies that with probability at least $1 - \text{poly}(\epsilon)$ the COMPARE call in Line 7 outputs either **High** or a value that is at least $\frac{7}{12} = \frac{1}{2} + \frac{1}{12}$. Since $\frac{7}{12} > \frac{1}{2} + 2^{j-5} \frac{\epsilon}{4}$ for $j = t$, the algorithm will output **REJECT** in Line 9.
- $j(H) < t$: By Equation (16) and the definition of the buckets, we have

$$\sum_{h \in H_{j(H)}} \left(\left(1 + 2^{j(H)} \frac{\epsilon}{4} \right) \frac{1}{N} - \frac{1}{N} \right) \geq \frac{\epsilon}{4t},$$

implying that $|H_{j(H)}| \geq \frac{N}{2^{j(H)} t}$ so that $D(H_{j(H)}) \geq \frac{1}{2^{j(H)} t}$. Therefore, when $j = j(H)$ so that $s_j = \Theta(2^{j(H)} t)$, with probability at least 99/100 the algorithm will get a point $h_{r'} \in H_{j(H)}$ in Line 4. Assume that indeed such a point $h_{r'}$ is selected. Since $D(h_{r'}) \geq (1 + 2^{j(H)-1} \frac{\epsilon}{4}) \frac{1}{N}$, while $D(i_r) \leq \frac{1}{N}$, for $\alpha_{j(H)} = 2^{j(H)-1} \frac{\epsilon}{4}$, we have

$$\frac{D(h_{r'})}{D(i_r)} \geq 1 + \alpha_{j(H)}.$$

Since COMPARE is called in Line 7 on the pair $\{i_r\}, \{h_{r'}\}$ with the “ δ ” parameter set to $\Theta(\epsilon 2^j)$, with probability $1 - \text{poly}(\epsilon)$ the algorithm outputs **REJECT** as a result of this COMPARE call.

Case 2: $|H| > \frac{N}{2}$. This proceeds similarly to Case 1. In this case we have that with high constant probability the algorithm selects a point $i_r \in H$ in Line 2. Here we consider the subcases $j(L) = t$ and $j(L) \leq t - 1$. In the first subcase we have that $\sum_{\ell \in L_t} \frac{1}{N} \geq \frac{\epsilon}{4t}$, so that $|L_t| \geq (\frac{\epsilon}{4t})n$, and in the second case we have that $\sum_{\ell \in L_{j(L)}} (2^{j(L)} \frac{\epsilon}{4}) \frac{1}{N} \geq \frac{\epsilon}{4t}$, so that $|L_{j(L)}| \geq \frac{N}{2^{j(L)}t}$. The analysis of each subcase is similar to Case 1. This concludes the proof of Theorem 4. \blacksquare

4.2 An $\Omega(1/\epsilon^2)$ lower bound for COND_D algorithms that test uniformity

In this subsection we give a lower bound showing that the query complexity of the PCOND_D algorithm of the previous subsection is essentially optimal, even for algorithms that may make general COND_D queries:

Theorem 5 *Any COND_D algorithm for testing whether $D = \mathcal{U}$ versus $d_{\text{TV}}(D, \mathcal{U}) \geq \epsilon$ must make $\Omega(1/\epsilon^2)$ queries.*

The high-level idea behind Theorem 5 is to reduce it to the well-known fact that distinguishing a fair coin from a $(\frac{1}{2} + 4\epsilon)$ -biased coin requires $\Omega(\frac{1}{\epsilon^2})$ coin tosses. We show that any q -query algorithm COND_D testing algorithm A can be transformed into an algorithm A' that successfully distinguishes q tosses of a fair coin from q tosses of a $(\frac{1}{2} + 4\epsilon)$ -biased coin.

Proof of Theorem 5: First note that we may assume without loss of generality that $0 < \epsilon \leq 1/8$. Let A be any q -query algorithm that makes COND_D queries and tests whether $D = \mathcal{U}$ versus $d_{\text{TV}}(D, \mathcal{U}) \geq \epsilon$. We may assume without loss of generality that in every possible execution algorithm A makes precisely q queries (this will be convenient later).

Let D_{No} be the distribution that has $D_{\text{No}}(i) = \frac{1+2\epsilon}{N}$ for each $i \in [1, \frac{N}{2}]$ and has $D_{\text{No}}(i) = \frac{1-2\epsilon}{N}$ for each $i \in [\frac{N}{2} + 1, N]$. (This is the “no”-distribution for our lower bound; it is ϵ -far in variation distance from the uniform distribution \mathcal{U} .) By Definition 2, it must be the case that

$$Z := \left| \Pr \left[A^{\text{COND}_{D_{\text{No}}}} \text{ outputs ACCEPT} \right] - \Pr \left[A^{\text{COND}_{\mathcal{U}}} \text{ outputs ACCEPT} \right] \right| \geq 1/3.$$

The proof works by showing that given A as described above, there must exist an algorithm A' with the following properties: A' is given as input a q -bit string $(b_1, \dots, b_q) \in \{0, 1\}^q$. Let D_0 denote the uniform distribution over $\{0, 1\}^q$ and let $D_{4\epsilon}$ denote the distribution over $\{0, 1\}^q$ in which each coordinate is independently set to 1 with probability $1/2 + 4\epsilon$. Then algorithm A' has

$$\left| \Pr_{b \sim D_0} [A'(b) \text{ outputs ACCEPT}] - \Pr_{b \sim D_{4\epsilon}} [A'(b) \text{ outputs ACCEPT}] \right| = Z. \quad (17)$$

Given (17), by the data processing inequality for total variation distance (Lemma 1) we have that $Z \leq d_{\text{TV}}(D_0, D_{4\epsilon})$. It is easy to see that $d_{\text{TV}}(D_0, D_{4\epsilon})$ is precisely equal to the variation distance $d_{\text{TV}}(\text{Bin}(q, 1/2), \text{Bin}(q, 1/2 + 4\epsilon))$. However, in order for the variation distance between these two binomial distributions to be as large as $1/3$ it must be the case that $q \geq \Omega(1/\epsilon^2)$:

Fact 4 (Distinguishing Fair from Biased Coin) Suppose $m \leq \frac{c}{\epsilon^2}$, with c a sufficiently small constant and $\epsilon \leq 1/8$. Then,

$$d_{\text{TV}}\left(\text{Bin}\left(m, \frac{1}{2}\right), \text{Bin}\left(m, \frac{1}{2} + 4\epsilon\right)\right) \leq \frac{1}{3}.$$

(Fact 4 is well known; it follows, for example, as an immediate consequence of Equations (2.15) and (2.16) of [AJ06].) Thus to prove Theorem 5 it remains only to describe algorithm A' and prove Equation (17).

As suggested above, algorithm A' uses algorithm A ; in order to do this, it must perfectly simulate the COND_D oracle that A requires, both in the case when $D = \mathcal{U}$ and in the case when $D = D_{\text{No}}$. We show below that when its input b is drawn from D_0 then A' can perfectly simulate the execution of A when it is run on the $\text{COND}_{\mathcal{U}}$ oracle, and when b is drawn from $D_{4\epsilon}$ then A' can perfectly simulate the execution of A when it is run on the $\text{COND}_{D_{\text{No}}}$ oracle.

Fix any step $1 \leq t \leq q$. We now describe how A' perfectly simulates the t -th step of the execution of A (i.e. the t -th call to COND_D that A makes, and the response of COND_D). We may inductively assume that A' has perfectly simulated the first $t - 1$ steps of the execution of A .

For each possible prefix of $t - 1$ query-response pairs to COND_D

$$\text{PREFIX} = ((S_1, s_1), \dots, (S_{t-1}, s_{t-1}))$$

(where each $S_i \subseteq [N]$ and each $s_i \in S_i$), there is some distribution $P_{A, \text{PREFIX}}$ over possible t -th query sets S_t that A would make given that its first $t - 1$ query-response pairs were PREFIX . So for a set $S \subseteq [n]$ and a possible prefix PREFIX , the value $P_{A, \text{PREFIX}}(S)$ is the probability that algorithm A , having had the transcript of its execution thus far be PREFIX , generates set S_t as its t -th query set. For any query set $S \subseteq [n]$, let us write S as a disjoint union $S = S_0 \amalg S_1$, where $S_0 = S \cap [1, \frac{N}{2}]$ and $S_1 = S \cap [\frac{N}{2} + 1, N]$. We may assume that every query S ever used by A has $|S_0|, |S_1| \geq 1$ (for otherwise A could perfectly simulate the response of $\text{COND}_D(S)$ whether D were \mathcal{U} or D_{No} by simply choosing a uniform point from S , so there would be no need to call COND_D on such an S). Thus we may assume that $P_{A, \text{PREFIX}}(S)$ is nonzero only for sets S that have $|S_0|, |S_1| \geq 1$.

Consider the bit $b_t \in \{0, 1\}$. As noted above, we inductively have that (whether D is \mathcal{U} or D_{No}) the algorithm A' has perfectly simulated the execution of A for its first $t - 1$ query-response pairs; in this simulation some prefix $\text{PREFIX} = ((S_1, s_1), \dots, (S_{t-1}, s_{t-1}))$ of query-response pairs has been constructed. If $b = (b_1, \dots, b_q)$ is distributed according to D_0 then PREFIX is distributed exactly according to the distribution of A 's prefixes of length $t - 1$ when A is run with $\text{COND}_{\mathcal{U}}$, and if $b = (b_1, \dots, b_q)$ is distributed according to $D_{4\epsilon}$ then the distribution of PREFIX is exactly the distribution of A 's prefixes of length $t - 1$ when A is run with $\text{COND}_{D_{\text{No}}}$.

Algorithm A' simulates the t -th stage of the execution of A as follows:

1. Randomly choose a set $S \subseteq [n]$ according to the distribution $P_{A, \text{PREFIX}}$; let $S = S_0 \amalg S_1$ be the set that is selected. Let us write $\alpha(S)$ to denote $|S_1|/|S_0|$ (so $\alpha(S) \in [2/N, N/2]$).
2. If $b_t = 1$ then set the bit $\sigma \in \{0, 1\}$ to be 1 with probability u_t and to be 0 with probability $1 - u_t$. If $b_t = 0$ then set σ to be 1 with probability v_t and to be 0 with probability $1 - v_t$. (We specify the exact values of u_t, v_t below.)

3. Set s to be a uniform random element of S_σ . Output the query-response pair $(S_t, s_t) = (S, s)$.

It is clear that Step 1 above perfectly simulates the t -th query that algorithm A would make (no matter what is the distribution D). To show that the t -th response is simulated perfectly, we must show that

- (i) if b_t is uniform random over $\{0, 1\}$ then s is distributed exactly as it would be distributed if A were being run on COND_U and had just proposed S as a query to COND_U ; i.e. we must show that s is a uniform random element of S_1 with probability $p(\alpha) \stackrel{\text{def}}{=} \frac{\alpha}{\alpha+1}$ and is a uniform random element of S_0 with probability $1 - p(\alpha)$.
- (ii) if $b_t \in \{0, 1\}$ has $\Pr[b_t = 1] = 1/2 + 4\epsilon$, then s is distributed exactly as it would be distributed if A were being run on $\text{COND}_{D_{N_0}}$ and had just proposed S as a query to COND_U ; i.e. we must show that s is a uniform random element of S_1 with probability $q(\alpha) \stackrel{\text{def}}{=} \frac{\alpha}{\alpha+(1+2\epsilon)/(1-2\epsilon)}$ and is a uniform random element of S_0 with probability $1 - q(\alpha)$.

By (i), we require that

$$\frac{u_t}{2} + \frac{v_t}{2} = p(\alpha) = \frac{\alpha}{\alpha+1}, \quad (18)$$

and by (ii) we require that

$$\left(\frac{1}{2} + 4\epsilon\right) u_t + \left(\frac{1}{2} - 4\epsilon\right) v_t = q(\alpha) = \frac{\alpha}{\alpha + \frac{1+2\epsilon}{1-2\epsilon}} \quad (19)$$

It is straightforward to check that

$$u_t = \frac{\alpha + 2\alpha^2 + 4\alpha\epsilon - 2\alpha^2\epsilon}{2 + 4\alpha + 2\alpha^2 + 4\epsilon - 2\alpha^2\epsilon}, \quad v_t = \frac{3\alpha + 2\alpha^2 + 4\alpha\epsilon - 2\alpha^2\epsilon}{2 + 4\alpha + 2\alpha^2 + 4\epsilon - 2\alpha^2\epsilon},$$

satisfy the above equations, and that for $0 < \alpha, 0 < \epsilon \leq 1/8$ we have $0 \leq u, v \leq 1$. So indeed A' perfectly simulates the execution of A in all stages $t = 1, \dots, q$. Finally, after simulating the t -th stage algorithm A' outputs whatever is output by its simulation of A , so Equation (17) indeed holds. This concludes the proof of Theorem 5. ■

5 Testing equivalence to a known distribution D^*

5.1 A $\text{poly}(\log n, 1/\epsilon)$ -query PCOND_D algorithm

In this subsection we present an algorithm PCOND-TEST-KNOWN and prove the following theorem:

Theorem 6 PCOND-TEST-KNOWN is a $\tilde{O}((\log N)^4/\epsilon^4)$ -query PCOND_D testing algorithm for testing equivalence to a known distribution D^* . That is, for every pair of distributions D, D^* over $[N]$ (such that D^* is fully specified and there is PCOND query access to D) the algorithm outputs ACCEPT with probability at least $2/3$ if $D = D^*$ and outputs REJECT with probability at least $2/3$ if $d_{\text{TV}}(D, D^*) \geq \epsilon$.

Intuition. Let D^* be a fully specified distribution, and let D be a distribution that may be accessed via a PCOND_D oracle. The high-level idea of the PCOND-TEST-KNOWN algorithm is the following: As in the case of testing uniformity, we shall try to “catch” a pair of points x, y such that $\frac{D(x)}{D(y)}$ differs significantly from $\frac{D^*(x)}{D^*(y)}$ (so that calling COMPARE_D on $\{x\}, \{y\}$ will reveal this difference). In the uniformity case, where $D^*(z) = 1/N$ for every z (so that $\frac{D^*(x)}{D^*(x)+D^*(y)} = 1/2$), to get a $\text{poly}(1/\epsilon)$ -query algorithm it was sufficient to show that sampling $\Theta(1/\epsilon)$ points uniformly (i.e., according to D^*) with high probability yields a point x for which $D(x) < D^*(x) - \Omega(\epsilon/N)$, and that sampling $\Theta(1/\epsilon)$ points from SAMP_D with high probability yields a point y for which $D(x) > D^*(y) + \Omega(\epsilon/N)$. However, for general D^* it is not sufficient to get such a pair because it is possible that $D^*(y)$ could be much larger than $D^*(x)$. If this were the case then it might happen that both $\frac{D^*(x)}{D^*(y)}$ and $\frac{D(x)}{D(y)}$ are very small, so calling COMPARE_D on $\{x\}, \{y\}$ cannot efficiently demonstrate that $\frac{D^*(x)}{D^*(y)}$ differs from $\frac{D(x)}{D(y)}$.

To address this issue we partition the points into $O(\log N/\epsilon)$ “buckets” so that within each bucket all points have similar probability according to D^* . We show that if D is ϵ -far from D^* , then either the probability weight of one of these buckets according to D differs significantly from what it is according to D^* (which can be observed by sampling from D), or we can get a pair $\{x, y\}$ that belong to the same bucket and for which $D(x)$ is sufficiently smaller than $D^*(x)$ and $D(y)$ is sufficiently larger than $D^*(y)$. For such a pair COMPARE will efficiently give evidence that D differs from D^* .

The algorithm and its analysis. We define some quantities that are used in the algorithm and its analysis. Let $\eta \stackrel{\text{def}}{=} \epsilon/c$ for some sufficiently large constant c that will be determined later. As described above we partition the domain elements $[N]$ into “buckets” according to their probability weight in D^* . Specifically, for $j = 1, \dots, \lceil \log(N/\eta) + 1 \rceil$, we let

$$B_j \stackrel{\text{def}}{=} \{x \in [N] : 2^{j-1} \cdot \eta/N \leq D^*(x) < 2^j \cdot \eta/N\}$$

and we let $B_0 \stackrel{\text{def}}{=} \{x \in [N] : D^*(x) < \eta/N\}$. Let $b \stackrel{\text{def}}{=} \lceil \log(N/\eta) + 1 \rceil + 1$ denote the number of buckets.

We further define $J^h \stackrel{\text{def}}{=} \{j : D^*(B_j) \geq \eta/b\}$ to denote the set of indices of “heavy” buckets, and let $J^\ell \stackrel{\text{def}}{=} \{j : D^*(B_j) < \eta/b\}$ denote the set of indices of “light” buckets. Note that we have

$$\sum_{j \in J^\ell \cup \{0\}} D^*(B_j) < 2\eta. \tag{20}$$

The query complexity of the algorithm is dominated by the number of PCOND_D queries performed in the executions of COMPARE , which by Lemma 2 is upper bounded by

$$O(s^2 \cdot b^2 \cdot (\log s)/\eta^2) = O\left(\frac{(\log \frac{N}{\epsilon})^4 \cdot \log((\log \frac{N}{\epsilon})/\epsilon)}{\epsilon^4}\right).$$

We argue completeness and soundness below.

Completeness: Suppose that $D = D^*$. Since the expected value of $\widehat{D}(B_j)$ (defined in Line 3) is precisely $D^*(B_j)$, for any fixed value of $j \in \{0, \dots, \lceil \log(N/\eta) + 1 \rceil\}$ an additive Chernoff bound

Algorithm 4: PCOND_D-TEST-KNOWN

Input: error parameter $\epsilon > 0$; query access to PCOND_D oracle; explicit description $(D^*(1), \dots, D^*(N))$ of distribution D^*

- 1: Call the SAMP_D oracle $m = \Theta(b^2(\log b)/\eta^2)$ times to obtain points h_1, \dots, h_m distributed according to D .
- 2: **for** $j = 0$ to b **do**
- 3: Let $\widehat{D}(B_j)$ be the fraction of points h_1, \dots, h_m that lie in B_j .
- 4: **if** some j has $|D^*(B_j) - \widehat{D}(B_j)| > \eta/b$ **then**
- 5: output REJECT and exit
- 6: **end if**
- 7: **end for**
- 8: Select $s = \Theta(b/\epsilon)$ points x_1, \dots, x_s independently from D^* .
- 9: Call the SAMP_D oracle $s = \Theta(b/\epsilon)$ times to obtain points y_1, \dots, y_s distributed according to D .
- 10: **for all** pairs (x_i, y_j) (where $1 \leq i, j \leq s$) such that $\frac{D^*(x)}{D^*(y)} \in [1/2, 2]$ **do**
- 11: Call COMPARE($\{x\}, \{y\}, \eta/(4b), 2, 1/(10s^2)$)
- 12: **if** COMPARE returns Low or a value smaller than $(1 - \eta/(2b)) \cdot \frac{D^*(x)}{D^*(y)}$ **then**
- 13: output REJECT (and exit)
- 14: **end if**
- 15: **end for**
- 16: output ACCEPT

implies that $|D^*(B_j) - \widehat{D}(B_j)| > \eta/b$ with failure probability at most $1/(10b)$. By a union bound over all b values of j , the algorithm outputs REJECT in Line 5 with probability at most $1/10$. Later in the algorithm, since $D = D^*$, no matter what points x_i, y_j are sampled from D^* and D respectively, the following holds for each pair (x_i, y_j) such that $D^*(x)/D^*(y) \in [1/2, 2]$. By Lemma 2 (and the setting of the parameters in the calls to COMPARE), the probability that COMPARE returns Low or a value smaller than $(1 - \delta/(2b)) \cdot (D^*(x)/D^*(y))$, is at most $1/(10s^2)$. A union bound over all (at most s^2) pairs (x_i, y_j) for which $D^*(x)/D^*(y) \in [1/2, 2]$, gives that the probability of outputting REJECT in Line 13 is at most $1/10$. Thus with overall probability at least $8/10$ the algorithm outputs ACCEPT.

Soundness: Now suppose that $d_{\text{TV}}(D, D^*) \geq \epsilon$; our goal is to show that the algorithm rejects with probability at least $2/3$. Since the algorithm rejects if any estimate $\widehat{D}(B_j)$ obtained in Line 3 deviates from $D^*(B_j)$ by more than $\pm\eta/b$, we may assume that all these estimates are indeed $\pm\eta/b$ -close to the values $D^*(B_j)$ as required. Moreover, by an additive Chernoff bound (as in the completeness analysis), we have that with overall failure probability at most $1/10$, each j has $|\widehat{D}(B_j) - D(B_j)| \leq \eta/b$; we condition on this event going forth. Thus, for every $0 \leq j \leq b$,

$$D^*(B_j) - 2\eta/b \leq D(B_j) \leq D^*(B_j) + 2\eta/b. \quad (21)$$

Recalling the definition of J^ℓ and Equation (20), we see that

$$\sum_{j \in J^\ell \cup \{0\}} D(B_j) < 4\eta. \quad (22)$$

Let

$$d_j \stackrel{\text{def}}{=} \sum_{x \in B_j} |D^*(x) - D(x)|, \quad (23)$$

so that $\|D^* - D\|_1 = \sum_j d_j$. By Equations (20) and (22), we have

$$\sum_{j \in J^\ell \cup \{0\}} d_j \leq \sum_{j \in J^\ell \cup \{0\}} (D^*(B_j) + D(B_j)) \leq 6\eta. \quad (24)$$

Since we have (by assumption) that $\|D^* - D\|_1 = 2d_{\text{TV}}(D^*, D) \geq 2\epsilon$, we get that

$$\sum_{j \in J^h \setminus \{0\}} d_j > 2\epsilon - 6\eta. \quad (25)$$

Let $N_j \stackrel{\text{def}}{=} |B_j|$ and observe that $N_j \leq D^*(B_j)/p_j \leq 1/p_j$, where $p_j \stackrel{\text{def}}{=} 2^{j-1} \cdot \eta/N$ is the lower bound on the probability (under D^*) of all elements in B_j . For each B_j such that $j \in J^h \setminus \{0\}$, let $H_j \stackrel{\text{def}}{=} \{x \in B_j : D(x) > D^*(x)\}$ and $L_j \stackrel{\text{def}}{=} \{x \in B_j : D(x) < D^*(x)\}$. Similarly to the “testing uniformity” analysis, we have that

$$\sum_{x \in L_j} (D^*(x) - D(x)) + \sum_{x \in H_j} (D(x) - D^*(x)) = d_j. \quad (26)$$

Equation (21) may be rewritten as

$$\left| \sum_{x \in L_j} (D^*(x) - D(x)) - \sum_{x \in H_j} (D(x) - D^*(x)) \right| \leq 2\eta/b, \quad (27)$$

and so we have both

$$\sum_{x \in L_j} (D^*(x) - D(x)) \geq d_j/2 - \eta/b \quad \text{and} \quad \sum_{x \in H_j} (D(x) - D^*(x)) \geq d_j/2 - \eta/b. \quad (28)$$

Also similarly to what we had before, let $H'_j \stackrel{\text{def}}{=} \{x \in B_j : D(x) > D^*(x) + \eta/(bN_j)\}$, and $L'_j \stackrel{\text{def}}{=} \{x \in B_j : D(x) < D^*(x) - \eta/(bN_j)\}$ (recall that $N_j = |B_j|$); these are the element of B_j that are “significantly heavier” (lighter, respectively) under D than under D^* . We have

$$\sum_{x \in L_j \setminus L'_j} (D^*(x) - D(x)) \leq \eta/b \quad \text{and} \quad \sum_{x \in H_j \setminus H'_j} (D(x) - D^*(x)) \leq \eta/b. \quad (29)$$

By Equation (25), there exists $j^* \in J^h \setminus \{0\}$ for which $d_{j^*} \geq (2\epsilon - 6\eta)/b$. For this index, applying Equations (28) and (29), we get that

$$\sum_{x \in L'_{j^*}} D^*(x) \geq \sum_{x \in L'_{j^*}} (D^*(x) - D(x)) \geq (\epsilon - 5\eta)/b, \quad (30)$$

and similarly,

$$\sum_{x \in H'_{j^*}} D(x) \geq \sum_{x \in H'_{j^*}} (D(x) - D^*(x)) \geq (\epsilon - 5\eta)/b. \quad (31)$$

Recalling that $\eta = \epsilon/6$, we have that $(\epsilon - 5\eta)/b = \epsilon/6b$. Now since $s = \Theta(b/\epsilon)$, with probability at least $9/10$ it is the case both that some x_i drawn in Line 8 belongs to L'_{j^*} and that some $y_{i'}$ drawn in Line 9 belongs to H'_{j^*} . By the definitions of L'_{j^*} and H'_{j^*} and the fact for each $j > 0$ it holds that $N_j \leq 1/p_j$ and $p_j \leq D^*(x) < 2p_j$ for each $x \in B_j$, we have that

$$D(x_i) < D^*(x_i) - \eta/(bN_{j^*}) \leq D^*(x_i) - (\eta/b)p_{j^*} \leq (1 - \eta/(2b))D^*(x_i) \quad (32)$$

and

$$D(y_{i'}) > D^*(y_{i'}) + \eta/(bN_{j^*}) \geq D^*(y_{i'}) + (\eta/b)p_j \geq (1 + \eta/(2b))D^*(y_{i'}). \quad (33)$$

Therefore,

$$\frac{D(x_i)}{D(y_{i'})} < \frac{1 - \eta/(2b)}{1 + \eta/(2b)} \cdot \frac{D^*(x_i)}{D^*(y_{i'})} < \left(1 - \frac{3\eta}{4b}\right) \cdot \frac{D^*(x_i)}{D^*(y_{i'})}. \quad (34)$$

By Lemma 2, with probability at least $1 - 1/(10s^2)$, the output of COMPARE is either **Low** or is at most $\left(1 - \frac{3\eta}{4b}\right) \cdot \left(1 + \frac{\eta}{4b}\right) < \left(1 - \frac{\eta}{2b}\right)$, causing the algorithm to reject. Thus the overall probability that the algorithm outputs **REJECT** is at least $8/10 - 1/(10s^2) > 2/3$, and the theorem is proved. \blacksquare

5.2 A $(\log N)^{\Omega(1)}$ lower bound for PCOND_D

In this subsection we prove that any PCOND_D algorithm for testing equivalence to a known distribution must have query complexity at least $(\log N)^{\Omega(1)}$:

Theorem 7 *Fix $\epsilon = 1/2$. There is a distribution D^* over $[N]$ (described below), which is such that any PCOND_D algorithm for testing whether $D = D^*$ versus $d_{\text{TV}}(D, D^*) \geq \epsilon$ must make $\Omega\left(\sqrt{\frac{\log N}{\log \log N}}\right)$ queries.*

The distribution D^* . Fix parameters $r = \Theta\left(\frac{\log N}{\log \log N}\right)$ and $K = \Theta(\log N)$. We partition $[N]$ from left (1) to right (N) into $2r$ consecutive intervals B_1, \dots, B_{2r} , which we henceforth refer to as “buckets.” The i -th bucket has $|B_i| = K^i$ (we may assume without loss of generality that N is of the form $\sum_{i=1}^r K^i$). The distribution D^* assigns equal probability weight to each bucket, so $D^*(B_i) = 1/(2r)$ for all $1 \leq i \leq 2r$. Moreover D^* is uniform within each bucket, so for all $j \in B_i$ we have $D^*(j) = 1/(2rK^i)$. This completes the specification of D^* .

To prove the lower bound we construct a probability distribution \mathcal{P}_{No} over possible “No”-distributions. To define the distribution \mathcal{P}_{No} it will be useful to have the notion of a “bucket-pair.” A bucket-pair U_i is $U_i = B_{2i-1} \cup B_{2i}$, i.e. the union of the i -th pair of consecutive buckets.

A distribution D drawn from \mathcal{P}_{No} is obtained by selecting a string $\pi = (\pi_1, \dots, \pi_r)$ uniformly at random from $\{\downarrow\uparrow, \uparrow\downarrow\}^r$ and setting D to be D_π , which we now define. The distribution D_π is obtained by perturbing D^* in the following way: for each bucket-pair $U_i = (B_{2i-1}, B_{2i})$,

- If $\pi_i = \uparrow\downarrow$ then the weight of B_{2i-1} is uniformly “scaled up” from $1/(2r)$ to $3/(4r)$ (keeping the distribution uniform within B_{2i-1}) and the weight of B_{2i} is uniformly “scaled down” from $1/(2r)$ to $1/(4r)$ (likewise keeping the distribution uniform within B_{2i}).
- If $\pi_i = \downarrow\uparrow$ then the weight of B_{2i-1} is uniformly “scaled down” from $1/(2r)$ to $1/(4r)$ and the weight of B_{2i} is uniformly “scaled up” from $1/(2r)$ to $3/(4r)$.

Note that for any distribution D in the support of \mathcal{P}_{No} and any $1 \leq i \leq r$ we have that $D(U_i) = D^*(U_i) = 1/r$.

Every distribution D in the support of \mathcal{P}_{No} has $d_{\text{TV}}(D^*, D) = 1/2$. Thus Theorem 7 follows immediately from the following:

Theorem 8 *Let A be any (possibly adaptive) algorithm. which makes at most $q \leq \frac{1}{3} \cdot \sqrt{r}$ calls to PCOND_D . Then*

$$\left| \Pr_{D \leftarrow \mathcal{P}_{\text{No}}} \left[A^{\text{PCOND}_D} \text{ outputs ACCEPT} \right] - \Pr \left[A^{\text{PCOND}_{D^*}} \text{ outputs ACCEPT} \right] \right| \leq 1/5. \quad (35)$$

Note that in the first probability of Equation (35) the randomness is over the draw of D from \mathcal{P}_{No} , the internal randomness of A in selecting its query sets, and the randomness of the responses to the PCOND_D queries. In the second probability the randomness is just over the internal coin tosses of A and the randomness of the responses to the PCOND_D queries.

Intuition for Theorem 8. A very high-level intuition for the lower bound is that PCOND_D queries are only useful for “comparing” points whose probabilities are within a reasonable multiplicative ratio of each other. But D^* and every distribution D in the support of \mathcal{P}_{No} are such that every two points either have the same probability mass under all of these distributions (so a PCOND_D query is not informative), or else the ratio of their probabilities is so skewed that a small number of PCOND_D queries is not useful for comparing them.

In more detail, we may suppose without loss of generality that in every possible execution, algorithm A first makes q calls to SAMP_D and then makes q (possibly adaptive) calls to PCOND_D . The more detailed intuition for the lower bound is as follows: First consider the SAMP_D calls. Since every possible D (whether D^* or a distribution drawn from \mathcal{P}_{No}) puts weight $1/r$ on each bucket-pair U_1, \dots, U_r , a birthday paradox argument implies that in both scenarios, with probability at least $9/10$ (over the randomness in the responses to the SAMP_D queries) no two of the $q \leq \frac{1}{3} \sqrt{r}$ calls to SAMP_D return points from the same bucket-pair. Conditioned on this, the distribution of responses to the SAMP_D queries is exactly the same under D^* and under D where D is drawn randomly from \mathcal{P}_{No} .

For the pair queries, the intuition is that in either setting (whether the distribution D is D^* or a randomly chosen distribution from \mathcal{P}_{No}), making q pair queries will with $1 - o(1)$ probability provide no information that the tester could not simulate for itself. This is because any pair query $\text{PCOND}_D(\{x, y\})$ either has x, y in the same bucket B_i or in different buckets $B_i \neq B_j$ with $i < j$. If x, y are both in the same bucket B_i then in either setting $\text{PCOND}_D(\{x, y\})$ is equally likely to return x or y . If they belong to buckets B_i, B_j with $i < j$ then in either setting $\text{PCOND}_D(\{x, y\})$ will return the one that belongs to P_i with probability $1 - 1/\Theta(K^{j-i}) \geq 1 - 1/\Omega(K)$.

Proof of Theorem 8: As described above, we may fix A to be any PCOND_D algorithm that makes exactly q calls to SAMP_D followed by exactly q adaptive calls to PCOND_D .

A *transcript* for A is a full specification of the sequence of interactions that A has with the PCOND_D oracle in a given execution. More precisely, it is a pair (Y, Z) where $Y = (s_1, \dots, s_q) \in [N]^q$ and $Z = ((\{x_1, y_1\}, p_1), \dots, (\{x_q, y_q\}, p_q))$, where $p_i \in \{x_i, y_i\}$ and $x_i, y_i \in [N]$. The idea is that Y is a possible sequence of responses that A might receive to the initial q SAMP_D queries, $\{x_i, y_i\}$ is a possible pair that could be the input to an i -th PCOND_D query, and p_i is a possible response that could be received from that query.

We say that a *length- i transcript prefix* is a pair (Y, Z^i) where Y is as above and $Z^i = ((\{x_1, y_1\}, p_1), \dots, (\{x_i, y_i\}, p_i))$. A PCOND algorithm A may be viewed as a collection of distributions over pairs $\{x, y\}$ in the following way: for each length- i transcript-prefix (Y, Z^i) ($0 \leq i \leq q-1$), there is a distribution over pairs $\{x_{i+1}, y_{i+1}\}$ that A would use to select the $(i+1)$ -st query pair for PCOND_D given that the length- i transcript prefix of A 's execution thus far was (Y, Z^i) . We write $T_{(Y, Z^i)}$ to denote this distribution over pairs.

Let P^* denote the distribution over transcripts induced by running A with oracle PCOND_{D^*} . Let P^{No} denote the distribution over transcripts induced by first (i) drawing D from \mathcal{P}_{No} , and then (ii) running A with oracle PCOND_D . To prove Theorem 8 it is sufficient to prove that the distribution over transcripts of A is statistically close whether the oracle is D^* or is a random D drawn from P^{No} , i.e. it is sufficient to prove that

$$d_{\text{TV}}(P^*, P^{\text{No}}) \leq 1/5. \quad (36)$$

For our analysis we will need to consider variants of algorithm A that, rather than making q calls to PCOND_D , instead “fake” the final $q - k$ of these PCOND_D queries as described below. For $0 \leq k \leq q$ we define $A^{(k)}$ to be the algorithm that works as follows:

1. $A^{(k)}$ exactly simulates the execution of A in making an initial q SAMP_D calls and making the first k PCOND_D queries precisely like A . Let (Y, Z^k) be the length- k transcript prefix of A 's execution thus obtained.
2. Exactly like A , algorithm $A^{(k)}$ draws a pair $\{x_{k+1}, y_{k+1}\}$ from $T_{(Y, Z^k)}$. However, instead of calling $\text{PCOND}_D(\{x_{k+1}, y_{k+1}\})$ to obtain p_{k+1} , algorithm $A^{(k)}$ generates p_{k+1} in the following manner:
 - (i) If x_{k+1} and y_{k+1} both belong to the same bucket B_ℓ then p_{k+1} is chosen uniformly from $\{x_{k+1}, y_{k+1}\}$.
 - (ii) If one of $\{x_{k+1}, y_{k+1}\}$ belongs to B_ℓ and the other belongs to $B_{\ell'}$ for some $\ell < \ell'$, then p_{k+1} is set to be the element of $\{x_{k+1}, y_{k+1}\}$ that belongs to B_ℓ .

Let (Y, Z^{k+1}) be the length- $(k+1)$ transcript prefix obtained by appending $(\{x_{k+1}, y_{k+1}\}, p_{k+1})$ to Z^k . Algorithm A' continues in this way for a total of $q - k$ stages; i.e. it next draws $\{x_{k+2}, y_{k+2}\}$ from $T_{(Y, Z^{k+1})}$ and generates p_{k+2} as described above; then (Y, Z^{k+2}) is the length- $(k+2)$ transcript prefix obtained by appending $(\{x_{k+2}, y_{k+2}\}, p_{k+2})$ to Z^{k+1} ; and so on. At the end of the process a transcript (Y, Z^q) has been constructed.

Let $P^{*,(k)}$ denote the distribution over final transcripts (Y, Z^q) that are obtained by running $A^{(k)}$ on a PCOND_{D^*} oracle. Let $P^{\text{No},(k)}$ denote the distribution over final transcripts (Y, Z^q) that are obtained by (i) first drawing D from \mathcal{P}_{No} , and then (ii) running $A^{(k)}$ on a PCOND_D oracle. Note that $P^{*,(q)}$ is identical to P^* and $P^{\text{No},(q)}$ is identical to P^{No} (since algorithm $A^{(q)}$, which does not fake any queries, is identical to algorithm A).

Recall that our goal is to prove Equation (36). Since $P^{*,(q)} = P^*$ and $P^{\text{No},(q)} = P^{\text{No}}$, Equation (36) is an immediate consequence (using the triangle inequality for total variation distance) of the following two lemmas, which we prove below:

Lemma 5 $d_{\text{TV}}(P^{*,(0)}, P^{\text{No},(0)}) \leq 1/10$.

Lemma 6 For all $0 \leq k < q$, we have $d_{\text{TV}}(P^{*,(k)}, P^{*,(k+1)}) \leq 1/(20q)$ and $d_{\text{TV}}(P^{\text{No},(k)}, P^{\text{No},(k+1)}) \leq 1/(20q)$.

Proof of Lemma 5: Define P_0^* to be the distribution over outcomes of the q calls to SAMP_D (i.e. over length-0 transcript prefixes) when $D = D^*$. Define P_0^{No} to be the distribution over outcomes of the q calls to SAMP_D when D is drawn from \mathcal{P}_{No} . We begin by noting that by the data processing inequality for total variation distance, we have $d_{\text{TV}}(P^{*,(0)}, P^{\text{No},(0)}) \leq d_{\text{TV}}(P_0^*, P_0^{\text{No}})$ (indeed, after the calls to respectively SAMP_D and SAMP_{D^*} , the same randomized function F – which fakes all remaining oracle calls – is applied to the two resulting distributions over length-0 transcript prefixes P_0^* and P_0^{No}). In the rest of the proof we show that $d_{\text{TV}}(P_0^*, P_0^{\text{No}}) \leq 1/10$.

Let E denote the event that the q calls to SAMP_D yield points s_1, \dots, s_q such that no bucket-pair U_i contains more than one of these points. Since $D^*(U_i) = 1/r$ for all i ,

$$P_0^*(E) = \prod_{j=0}^{q-1} \left(1 - \frac{j}{r}\right) \geq 9/10, \quad (37)$$

where Equation (37) follows from a standard birthday paradox analysis and the fact that $q \leq \frac{1}{3}\sqrt{r}$. Since for each possible outcome of D drawn from \mathcal{P}_{No} we have $D(U_i) = 1/r$ for all i , we further have that also

$$P_0^{\text{No}}(E) = \prod_{j=0}^{q-1} \left(1 - \frac{j}{r}\right). \quad (38)$$

We moreover claim that the two conditional distributions $(P_0^*|E)$ and $(P_0^{\text{No}}|E)$ are identical, i.e.

$$(P_0^*|E) = (P_0^{\text{No}}|E). \quad (39)$$

To see this, fix any sequence $(\ell_1, \dots, \ell_q) \in [r]^q$ such that $\ell_i \neq \ell_j$ for all $i \neq j$. Let $(s_1, \dots, s_q) \in [N]^q$ denote a draw from $(P_0^*|E)$. The probability that $(s_i \in U_{\ell_i}$ for all $1 \leq i \leq q)$ is precisely $1/r^q$. Now given that $s_i \in U_{\ell_i}$ for all i , it is clear that s_i is equally likely to lie in $B_{2\ell_i-1}$ and in $B_{2\ell_i}$, and given that it lies in a particular one of the two buckets, it is equally likely to be any element in that bucket. This is true independently for all $1 \leq i \leq q$.

Now let $(s_1, \dots, s_q) \in [N]^q$ denote a draw from $(P_0^{\text{No}}|E)$. Since each distribution D in the support of \mathcal{P}_{No} has $D(U_i) = 1/r$ for all i , we likewise have that the probability that $(s_i \in U_{\ell_i}$ for all $1 \leq i \leq q)$

is precisely $1/r^q$. Now given that $s_i \in U_{\ell_i}$ for all i , we have that s_i is equally likely to lie in $B_{2\ell_i-1}$ and in $B_{2\ell_i}$; this is because π_i (recall that π determines $D = D_\pi$) is equally likely to be $\uparrow\downarrow$ (in which case $D(B_{2\ell_i-1}) = 3/(4r)$ and $D(B_{2\ell_i}) = 1/(4r)$) as it is to be $\downarrow\uparrow$ (in which case $D(B_{2\ell_i-1}) = 1/(4r)$ and $D(B_{2\ell_i}) = 3/(4r)$). Additionally, given that s_i lies in a particular one of the two buckets, it is equally likely to be any element in that bucket. This is true independently for all $1 \leq i \leq q$ (because conditioning on E ensures that no two elements of s_1, \dots, s_q lie in the same bucket-pair, so there is “fresh randomness for each i ”), and so indeed the two conditional distributions $(P_0^*|E)$ and $(P_0^{\text{No}}|E)$ are identical.

Finally, the claimed bound $d_{\text{TV}}(P_0^*, P_0^{\text{No}}) \leq 1/10$ follows directly from Equations (37), (38) and (39). ■

Proof of Lemma 6: Consider first the claim that $d_{\text{TV}}(P^{*,(k)}, P^{*,(k+1)}) \leq 1/(20q)$. Fix any $0 \leq k < q$. The data processing inequality for total variation distance implies that $d_{\text{TV}}(P^{*,(k)}, P^{*,(k+1)})$ is at most the variation distance between random variables X and Y , where

- X is the random variable obtained by running A on COND_{D^*} to obtain a length- k transcript prefix (Y, Z^k) , then drawing $\{x_{k+1}, y_{k+1}\}$ from $\mathbb{T}_{(Y, Z^k)}$, then setting p_{k+1} to be the output of $\text{PCOND}_{D^*}(\{x_{k+1}, y_{k+1}\})$; and
- Y is the random variable obtained by running A on COND_{D^*} to obtain a length- k transcript prefix (Y, Z^k) , then drawing $\{x_{k+1}, y_{k+1}\}$ from $\mathbb{T}_{(Y, Z^k)}$, then setting p_{k+1} according to the rules 2(i) and 2(ii) given above.

Consider any fixed outcome of (Y, Z^k) and $\{x_{k+1}, y_{k+1}\}$. If rule 2(i) is applied (x_{k+1} and y_{k+1} are in the same bucket) then there is zero contribution to the variation distance between X and Y , because choosing a uniform element of $\{x_{k+1}, y_{k+1}\}$ is a perfect simulation of $\text{PCOND}_D(\{x_{k+1}, y_{k+1}\})$. If rule 2(ii) is applied then the contribution is at most $O(1/K) < 1/20q$, because $\text{PCOND}_{D^*}(\{x_{k+1}, y_{k+1}\})$ would return a different outcome from rule 2(ii) with probability $1/\Theta(K^{\ell'-\ell}) = O(1/K)$. Averaging over all possible outcomes of (Y, Z^k) and $\{x_{k+1}, y_{k+1}\}$ we get that the variation distance between X and Y is at most $1/20q$ as claimed.

An identical argument shows that similarly $d_{\text{TV}}(P^{\text{No},(k)}, P^{\text{No},(k+1)}) \leq 1/(20q)$. The key observation is that for any distribution D in the support of P^{No} , as with D^* it is the case that points in the same bucket have equal probability under D and a point y that is $\ell' - \ell$ buckets lower than x has probability only $1/\Theta(K^{\ell'-\ell})$ of being returned by a call to $\text{PCOND}_D(\{x, y\})$. This concludes the proof of Lemma 6 and of Theorem 7. ■

5.3 A $\text{poly}(1/\epsilon)$ -query COND_D algorithm

In this subsection we present an algorithm COND-TEST-KNOWN and prove the following theorem:

Theorem 9 COND-TEST-KNOWN is a $\tilde{O}(1/\epsilon^4)$ -query COND_D testing algorithm for testing equivalence to a known distribution D^* . That is, for every pair of distributions D, D^* over $[N]$ (such that D^* is fully specified and there is COND query access to D), the algorithm outputs ACCEPT with probability at least $2/3$ if $D = D^*$ and outputs REJECT with probability at least $2/3$ if $d_{\text{TV}}(D, D^*) \geq \epsilon$.

This constant-query testing algorithm stands in interesting contrast to the $(\log N)^{\Omega(1)}$ -query lower bound for PCOND_D algorithms for this problem.

High-level overview of the algorithm and its analysis: First, we note that by reordering elements of $[N]$ we may assume without loss of generality that $D^*(1) \leq \dots \leq D^*(N)$; this will be convenient for us.

Our $(\log N)^{\Omega(1)}$ query lower bound for PCOND_D algorithms exploited the intuition that comparing two points using the PCOND_D oracle might not provide much information (e.g. if one of the two points was a priori “known” to be much heavier than the other). In contrast, with a general COND_D oracle at our disposal, we can compare a given point $j \in [N]$ with *any subset* of $[N] \setminus \{j\}$. Thus the following definition will be useful:

Definition 3 (comparable points) Fix $0 < \lambda \leq 1$. A point $j \in \text{supp}(D^*)$ is said to be λ -comparable if there exists a set $S \subseteq ([N] \setminus \{j\})$ such that

$$D^*(j) \in [\lambda D^*(S), D^*(S)/\lambda].$$

Such a set S is then said to be a λ -comparable-witness for j (according to D^*), which is denoted $S \cong^* j$. We say that a set $T \subseteq [N]$ is λ -comparable if every $i \in T$ is λ -comparable.

We stress that the notion of being λ -comparable deals only with the known distribution D^* ; this will be important later.

Fix $\epsilon_1 = \Theta(\epsilon)$ (we specify ϵ_1 precisely in Equation 42 below). Our analysis and algorithm consider two possible cases for the distribution D^* (where it is not hard to verify, and we provide an explanation subsequently, that one of the two cases must hold):

1. The first case is that for some $i^* \in [N]$ we have

$$D^*({1, \dots, i^*}) > 2\epsilon_1 \quad \text{but} \quad D^*({1, \dots, i^* - 1}) \leq \epsilon_1. \quad (40)$$

In this case $1 - \epsilon_1$ of the total probability mass of D^* must lie on a set of at most $1/\epsilon_1$ elements, and in such a situation it is easy to efficiently test whether $D = D^*$ using $\text{poly}(1/\epsilon)$ queries (see Algorithm $\text{COND}_D\text{-TEST-KNOWN-HEAVY}$ and Lemma 9).

2. The second case is that there exists an element $k^* \in [N]$ such that

$$\epsilon_1 < D^*({1, \dots, k^*}) \leq 2\epsilon_1 < D^*({1, \dots, k^* + 1}). \quad (41)$$

This is the more challenging (and typical) case. In this case, it can be shown that every element $j > k^*$ has at least one ϵ_1 -comparable-witness within $\{1, \dots, j\}$. In fact, we show (see Claim 7) that either (a) $\{1, \dots, j - 1\}$ is an ϵ_1 -comparable witness for j , or (b) the set $\{1, \dots, j - 1\}$ can be partitioned into disjoint sets⁹ S_1, \dots, S_t such that *each* S_i , $1 \leq i \leq t$, is a $\frac{1}{2}$ -comparable-witness for j . Case (a) is relatively easy to handle so we focus on (b) in our informal description below.

⁹In fact the sets are intervals (under the assumption $D^*(1) \leq \dots \leq D^*(n)$), but that is not really important for our arguments.

The partition S_1, \dots, S_t is useful to us for the following reason: Suppose that $d_{\text{TV}}(D, D^*) \geq \epsilon$. It is not difficult to show (see Claim 8) that unless $D(\{1, \dots, k^*\}) > 3\epsilon_1$ (which can be easily detected and provides evidence that the tester should reject), a random sample of $\Theta(1/\epsilon)$ draws from D will with high probability contain a “heavy” point $j > k^*$, that is, a point $j > k^*$ such that $D(j) \geq (1 + \epsilon_2)D^*(j)$ (where $\epsilon_2 = \Theta(\epsilon)$). Given such a point j , there are two possibilities:

1. The first possibility is that a significant fraction of the sets S_1, \dots, S_t have $D(j)/D(S_i)$ “noticeably different” from $D^*(j)/D^*(S_i)$. (Observe that since each set S_i is a $\frac{1}{2}$ -comparable witness for j , it is possible to efficiently check whether this is the case.) If this is the case then our tester should reject since this is evidence that $D \neq D^*$.
2. The second possibility is that almost every S_i has $D(j)/D(S_i)$ very close to $D^*(j)/D^*(S_i)$. If this is the case, though, then since $D(j) \geq (1 + \epsilon_2)D^*(j)$ and the union of S_1, \dots, S_t is $\{1, \dots, j-1\}$, it must be the case that $D(\{1, \dots, j\})$ is “significantly larger” than $D^*(\{1, \dots, j\})$. This will be revealed by random sampling from D and thus our testing algorithm can reject in this case as well.

Key quantities and useful claims. We define some quantities that are used in the algorithm and its analysis. Let

$$\epsilon_1 \stackrel{\text{def}}{=} \frac{\epsilon}{10}; \quad \epsilon_2 \stackrel{\text{def}}{=} \frac{\epsilon}{2}; \quad \epsilon_3 \stackrel{\text{def}}{=} \frac{\epsilon}{48}; \quad \epsilon_4 \stackrel{\text{def}}{=} \frac{\epsilon}{6}. \quad (42)$$

Claim 7 *Suppose there exists an element $k^* \in [N]$ that satisfies Equation (41). Fix any $j > k^*$. Then*

1. *If $D^*(j) \geq \epsilon_1$, then $S_1 \stackrel{\text{def}}{=} \{1, \dots, j-1\}$ is an ϵ_1 -comparable witness for j ;*
2. *If $D^*(j) < \epsilon_1$ then the set $\{1, \dots, j-1\}$ can be partitioned into disjoint sets S_1, \dots, S_t such that each S_i , $1 \leq i \leq t$, is a $\frac{1}{2}$ -comparable-witness for j .*

Proof: First consider the case that $D^*(j) \geq \epsilon_1$. In this case $S_1 = \{1, \dots, j-1\}$ is an ϵ_1 -comparable witness for j because $D^*(j) \geq \epsilon_1 \geq \epsilon_1 D^*(\{1, \dots, j-1\})$ and $D^*(j) \leq 1 \leq \frac{1}{\epsilon_1} D^*(\{1, \dots, k^*\}) \leq \frac{1}{\epsilon_1} D^*(\{1, \dots, j-1\})$, where the last inequality holds since $k^* \leq j-1$.

Next, consider the case that $D^*(j) < \epsilon_1$. In this case we build our intervals iteratively from right to left, as follows. Let $j_1 = j-1$ and let j_2 be the minimum index in $\{0, \dots, j_1-1\}$ such that

$$D^*(\{j_2+1, \dots, j_1\}) \leq D^*(j).$$

(Observe that we must have $j_2 \geq 1$, because $D^*(\{1, \dots, k^*\}) > \epsilon_1 > D^*(j)$.) Since $D^*(\{j_2, \dots, j_1\}) > D^*(j)$ and the function $D^*(\cdot)$ is monotonically increasing, it must be the case that

$$\frac{1}{2}D^*(j) \leq D^*(\{j_2+1, \dots, j_1\}) \leq D^*(j).$$

Thus the interval $S_1 \stackrel{\text{def}}{=} \{j_2+1, \dots, j_1\}$ is a $\frac{1}{2}$ -comparable witness for j as desired.

We continue in this fashion from right to left; i.e. if we have defined j_2, \dots, j_t as above and there is an index $j' \in \{0, \dots, j_t - 1\}$ such that $D^*({j' + 1, \dots, j_t}) > D^*(j)$, then we define j_{t+1} to be the minimum index in $\{0, \dots, j_t - 1\}$ such that

$$D^*({j_{t+1} + 1, \dots, j_t}) \leq D^*(j),$$

and we define S_t to be the interval $\{j_{t+1} + 1, \dots, j_t\}$. The argument of the previous paragraph tells us that

$$\frac{1}{2}D^*(j) \leq D^*({j_{t+1} + 1, \dots, j_t}) \leq D^*(j) \quad (43)$$

and hence S_t is an $\frac{1}{2}$ -comparable witness for j .

At some point, after intervals $S_1 = \{j_2 + 1, \dots, j_1\}, \dots, S_t = \{j_{t+1} + 1, \dots, j_t\}$ have been defined in this way, it will be the case that there is no index $j' \in \{0, \dots, j_t - 1\}$ such that $D^*({j' + 1, \dots, j_t}) > D^*(j)$. At this point there are two possibilities: first, if $j_{t+1} + 1 = 1$, then S_1, \dots, S_t give the desired partition of $\{1, \dots, j - 1\}$. If $j_{t+1} + 1 > 1$ then it must be the case that $D^*({1, \dots, j_{t+1}}) \leq D^*(j)$. In this case we simply add the elements $\{1, \dots, j_{t+1}\}$ to S_t , i.e. we redefine S_t to be $\{1, \dots, j_t\}$. By Equation (43) we have that

$$\frac{1}{2}D^*(j) \leq D^*(S_t) \leq 2D^*(j)$$

and thus S_t is an $\frac{1}{2}$ -comparable witness for j as desired. This concludes the proof. \blacksquare

Definition 4 (Heavy points) A point $j \in \text{supp}(D^*)$ is said to be η -heavy if $D(j) \geq (1 + \eta)D^*(j)$.

Claim 8 Suppose that $d_{\text{TV}}(D, D^*) \geq \epsilon$ and Equation (41) holds. Suppose moreover that $D(\{1, \dots, k^*\}) \leq 4\epsilon_1$. Let i_1, \dots, i_ℓ be i.i.d. points drawn from D . Then for $\ell = \Theta(1/\epsilon)$, with probability at least 99/100 (over the i.i.d. draws of $i_1, \dots, i_\ell \sim D$) there is some point $i_j \in \{i_1, \dots, i_\ell\}$ such that $i_j > k^*$ and i_j is ϵ_2 -heavy.

Proof: Define H_1 to be the set of all ϵ_2 -heavy points and H_2 to be the set of all “slightly lighter” points as follows:

$$\begin{aligned} H_1 &= \{ i \in [N] \mid D(i) \geq (1 + \epsilon_2)D^*(i) \} \\ H_2 &= \{ i \in [N] \mid (1 + \epsilon_2)D^*(i) > D(i) \geq D^*(i) \} \end{aligned}$$

By definition of the total variation distance, we have

$$\begin{aligned} \epsilon \leq d_{\text{TV}}(D, D^*) &= \sum_{i: D(i) \geq D^*(i)} (D(i) - D^*(i)) = (D(H_1) - D^*(H_1)) + (D(H_2) - D^*(H_2)) \\ &\leq D(H_1) + ((1 + \epsilon_2)D^*(H_2) - D^*(H_2)) \\ &= D(H_1) + \epsilon_2 D^*(H_2) < D(H_1) + \epsilon_2 = D(H_1) + \frac{\epsilon}{2}. \end{aligned}$$

So it must be the case that $D(H_1) \geq \epsilon/2 = 5\epsilon_1$. Since by assumption we have $D(\{1, \dots, k^*\}) \leq 4\epsilon_1$, it must be the case that $D(H_1 \setminus \{1, \dots, k^*\}) \geq \epsilon_1$. The claim follows from the definition of H_1 and the size, ℓ , of the sample. \blacksquare

Algorithm 5: COND_D-TEST-KNOWN

Input: error parameter $\epsilon > 0$; query access to COND_D oracle; explicit description $(D^*(1), \dots, D^*(N))$ of distribution D^* satisfying $D^*(1) \leq \dots \leq D^*(N)$

- 1: Let i^* be the minimum index $i \in [N]$ such that $D^*({1, \dots, i}) > 2\epsilon_1$.
- 2: **if** $D^*({1, \dots, i^* - 1}) \leq \epsilon_1$ **then**
- 3: Call algorithm COND_D-Test-Known-Heavy(ϵ , COND_D, D^* , i^*) (and exit)
- 4: **else**
- 5: Call algorithm COND_D-Test-Known-Main(ϵ , COND_D, D^* , $i^* - 1$) (and exit).
- 6: **end if**

Algorithm 6: COND_D-TEST-KNOWN-HEAVY

Input: error parameter $\epsilon > 0$; query access to COND_D oracle; explicit description $(D^*(1), \dots, D^*(N))$ of distribution D^* satisfying $D^*(1) \leq \dots \leq D^*(N)$; value $i^* \in [N]$ satisfying $D^*({1, \dots, i^* - 1}) \leq \epsilon_1$, $D^*({1, \dots, i^*}) > 2\epsilon_1$

- 1: Call the SAMP_D oracle $m = \Theta((\log(1/\epsilon))/\epsilon^4)$ times. For each $i \in [i^*, N]$ let $\widehat{D}(j)$ be the fraction of the m calls to SAMP_D that returned i . Let $\widehat{D}' = 1 - \sum_{i \in [i^*, N]} \widehat{D}(i)$ be the fraction of the m calls that returned values in $\{1, \dots, i^* - 1\}$.
- 2: **if** either (any $i \in [i^*, N]$ has $|\widehat{D}(i) - D^*(i)| > \epsilon_1^2$) or $(\widehat{D}' - D^*({1, \dots, i^* - 1})) > \epsilon_1$ **then**
- 3: output REJECT (and exit)
- 4: **end if**
- 5: Output ACCEPT

5.3.1 Proof of Theorem 6

It is straightforward to verify that the query complexity of COND_D-Test-Known-Heavy is $\tilde{O}(1/\epsilon^4)$ and the query complexity of COND_D-Test-Known-Main is also $\tilde{O}(1/\epsilon^4)$, so the overall query complexity of COND-TEST-KNOWN is as claimed.

By the definition of i^* (in the first line of the algorithm), either Equation (40) holds for this setting of i^* , or Equation (41) holds for $k^* = i^* - 1$. To prove correctness of the algorithm, we first deal with the simpler case, which is that Equation (40) holds:

Lemma 9 *Suppose that D^* is such that $D^*({1, \dots, i^*}) > 2\epsilon_1$ but $D^*({1, \dots, i^* - 1}) \leq \epsilon_1$. Then COND_D-TEST-KNOWN-HEAVY(ϵ , COND_D, D^* , i^*) returns ACCEPT with probability at least 2/3 if $D = D^*$ and returns REJECT with probability at least 2/3 if $d_{\text{TV}}(D, D^*) \geq \epsilon$.*

Proof: The conditions of Lemma 9, together with the fact that $D^*(\cdot)$ is monotone non-decreasing, imply that each $i \geq i^*$ has $D^*(i) \geq \epsilon_1$. Thus there can be at most $1/\epsilon_1$ many values $i \in \{i^*, \dots, N\}$, i.e. it must be the case that $i^* \geq N - 1/\epsilon_1 + 1$. Since the expected value of $\widehat{D}(i)$ (defined in Line 1 of COND_D-TEST-KNOWN-HEAVY) is precisely $D(i)$, for any fixed value of $i \in \{i^*, \dots, n\}$ an additive Chernoff bound implies that $|D(i) - \widehat{D}(i)| \leq (\epsilon_1)^2$ with failure probability at most $\frac{1}{10(1+\frac{1}{\epsilon_1})}$. Similarly $|\widehat{D}' - D^*({1, \dots, i^* - 1})| \leq \epsilon_1$ with failure probability at most $\frac{1}{10(1+\frac{1}{\epsilon_1})}$. A union bound over all failure events gives that with probability at least 9/10 each value $i \in \{i^*, \dots, N\}$

Algorithm 7: COND_D-TEST-KNOWN-MAIN

- Input:** error parameter $\epsilon > 0$; query access to COND_D oracle; explicit description $(D^*(1), \dots, D^*(N))$ of distribution D^* satisfying $D^*(1) \leq \dots \leq D^*(N)$; value $k^* \in [N]$ satisfying $\epsilon_1 < D^*({1, \dots, k^*}) \leq 2\epsilon_1 < D^*({1, \dots, k^* + 1})$
- 1: Call the SAMP_D oracle $\Theta(1/\epsilon^2)$ times and let $\widehat{D}({1, \dots, k^*})$ denote the fraction of responses that lie in ${1, \dots, k^*}$. If $\widehat{D}({1, \dots, k^*}) \notin [\frac{\epsilon_1}{2}, \frac{5\epsilon_1}{2}]$ then output REJECT (and exit).
 - 2: Call the SAMP_D oracle $\ell = \Theta(1/\epsilon)$ times to obtain points i_1, \dots, i_ℓ .
 - 3: **for** all $j \in {1, \dots, \ell}$ such that $i_j > k^*$ **do**
 - 4: Call the SAMP_D oracle $m = \Theta(\log(1/\epsilon)/\epsilon)$ times and let $\widehat{D}({1, \dots, i_j})$ be the fraction of responses that lie in ${1, \dots, i_j}$. If $\widehat{D}({1, \dots, i_j}) \notin [1 - \epsilon_3, 1 + \epsilon_3]D^*({1, \dots, i_j})$ then output REJECT (and exit).
 - 5: **if** $D^*(i_j) \geq \epsilon_1$ **then**
 - 6: Run COMPARE(${i_j}, {1, \dots, i_j - 1}, \frac{\epsilon_2}{16}, \frac{2}{\epsilon_1}, \frac{1}{10\ell}$) and let v denote its output. If $v \notin [1 - \frac{\epsilon_2}{8}, 1 + \frac{\epsilon_2}{8}] \frac{D^*({1, \dots, i_j - 1})}{D^*({i_j})}$ then output REJECT (and exit).
 - 7: **else**
 - 8: Let S_1, \dots, S_t be the partition of ${1, \dots, i_j - 1}$ such that each S_i is an ϵ_1 -comparable witness for i_j , which is provided by Claim 7.
 - 9: Select a list of $h = \Theta(1/\epsilon)$ elements S_{a_1}, \dots, S_{a_h} independently and uniformly from ${S_1, \dots, S_t}$.
 - 10: For each $S_{a_r}, 1 \leq r \leq h$, run COMPARE(${i_j}, S_{a_r}, \frac{\epsilon_4}{8}, 4, \frac{1}{10\ell h}$) and let v denote its output. If $v \notin [1 - \frac{\epsilon_4}{4}, 1 + \frac{\epsilon_4}{4}] \frac{D^*(S_{a_r})}{D^*({i_j})}$ then output REJECT (and exit).
 - 11: **end if**
 - 12: **end for**
 - 13: Output ACCEPT.
-

has $|D(i) - \widehat{D}(i)| \leq \epsilon_1^2$ and additionally $|\widehat{D}' - D(\{1, \dots, i^* - 1\})| \leq \epsilon_1$; we refer to this compound event as (*).

If $D^* = D$, by (*) the algorithm outputs ACCEPT with probability at least 9/10.

Now suppose that $d_{\text{TV}}(D, D^*) \geq \epsilon$. With probability at least 9/10 we have (*) so we suppose that indeed (*) holds. In this case we have

$$\begin{aligned}
\epsilon \leq d_{\text{TV}}(D, D^*) &= \sum_{i < i^*} |D(i) - D^*(i)| + \sum_{i \geq i^*} |D(i) - D^*(i)| \\
&\leq \sum_{i < i^*} (D(i) + D^*(i)) + \sum_{i \geq i^*} |D(i) - D^*(i)| \\
&\leq D(\{1, \dots, i^* - 1\}) + \epsilon_1 + \sum_{i \geq i^*} (|\widehat{D}(i) - D^*(i)| + \epsilon_1^2) \\
&\leq \widehat{D}' + \epsilon_1 + 2\epsilon_1 + \sum_{i \geq i^*} (|\widehat{D}(i) - D^*(i)|)
\end{aligned}$$

where the first inequality is by the triangle inequality, the second is by (*) and the fact that $D^*(\{1, \dots, i^* - 1\}) \leq \epsilon_1$, and the third inequality is by (*) and the fact that there are at most $1/\epsilon_1$ elements in $\{i^*, \dots, N\}$. Since $\epsilon_1 = \epsilon/10$, the above inequality implies that

$$\frac{7}{10}\epsilon \leq \widehat{D}' + \sum_{i \geq i^*} (|\widehat{D}(i) - D^*(i)|).$$

If any $i \in \{i^*, \dots, N\}$ has $|\widehat{D}(i) - D^*(i)| > (\epsilon_1)^2$ then the algorithm outputs REJECT so we may assume that $|\widehat{D}(i) - D^*(i)| \leq \epsilon_1^2$ for all i . This implies that

$$6\epsilon_1 = \frac{6}{10}\epsilon \leq \widehat{D}'$$

but since $D^*(\{1, \dots, i^* - 1\}) \leq \epsilon_1$ the algorithm must REJECT. ■

Now we turn to the more difficult (and typical) case, that Equation (41) holds (for $k^* = i^* - 1$), i.e.

$$\epsilon_1 < D^*(\{1, \dots, k^*\}) \leq 2\epsilon_1 < D^*(\{1, \dots, k^* + 1\}).$$

With the claims we have already established it is straightforward to argue completeness:

Lemma 10 *Suppose that $D = D^*$ and Equation (41) holds. Then with probability at least 2/3 algorithm COND_D-TEST-KNOWN-MAIN outputs ACCEPT.*

Proof: We first observe that the expected value of the quantity $\widehat{D}(\{1, \dots, k^*\})$ defined in Line 1 is precisely $D(\{1, \dots, k^*\}) = D^*(\{1, \dots, k^*\})$ and hence lies in $[\epsilon_1, 2\epsilon_1]$ by Equation (41). The additive Chernoff bound implies that the probability the algorithm outputs REJECT in Line 1 is at most 1/10. Thus we may assume the algorithm continues to Line 2.

In any given execution of Line 4, since the expected value of $\widehat{D}(\{1, \dots, i_j\})$ is precisely $D(\{1, \dots, i_j\}) = D^*(\{1, \dots, i_j\}) > \epsilon_1$, a multiplicative Chernoff bound gives that the algorithm

outputs REJECT with probability at most $1/(10\ell)$. Thus the probability that the algorithm outputs REJECT in any execution of Line 4 is at most $1/10$. We henceforth assume that the algorithm never outputs REJECT in this step.

Fix a setting of $j \in \{1, \dots, \ell\}$ such that $i_j > k^*$. Consider first the case that $D^*(i_j) \geq \epsilon_1$ so the algorithm enters Line 6. By item (1) of Claim 7 and item (1) of Lemma 2, we have that with probability at least $1 - \frac{1}{10\ell}$ COMPARE outputs a value v in the range $[1 - \frac{\epsilon_2}{16}, 1 + \frac{\epsilon_2}{16}] \frac{D^*(\{1, \dots, i_j-1\})}{D^*(\{i_j\})}$ (recall that $D = D^*$), so the algorithm does not output REJECT in Line 6. Now suppose that $D^*(i_j) < \epsilon_1$ so the algorithm enters Line 8. Fix a value $1 \leq r \leq h$ in Line 10. By Claim 7 we have that S_{a_r} is a $\frac{1}{2}$ -comparable witness for i_j . By item (1) of Lemma 2, we have that with probability at least $1 - \frac{1}{10\ell h}$ COMPARE outputs a value v in the range $[1 - \frac{\epsilon_4}{4}, 1 + \frac{\epsilon_4}{4}] \frac{D^*(S_{a_r})}{D^*(\{i_j\})}$ (recall that $D = D^*$). A union bound over all h values of r gives that the algorithm outputs REJECT in Line 10 with probability at most $1/(10\ell)$. So in either case, for this setting of j , the algorithm outputs REJECT on that iteration of the outer loop with probability at most $1/(10\ell)$. A union bound over all ℓ iterations of the outer loop gives that the algorithm outputs REJECT at any execution of Line 6 or Line 10 is at most $1/10$.

Thus the overall probability that the algorithm outputs REJECT is at most $3/10$, and the lemma is proved. ■

Next we argue soundness:

Lemma 11 *Suppose that $d_{\text{TV}}(D, D^*) \geq \epsilon$ and Equation (41) holds. Then with probability at least $2/3$ algorithm COND_D-TEST-KNOWN-MAIN outputs REJECT.*

Proof: If $D(\{1, \dots, k^*\}) \notin [\epsilon_1, 3\epsilon_1]$ then a standard additive Chernoff bound implies that the algorithm outputs REJECT in Line 1 with probability at least $9/10$. Thus we may assume going forward in the argument that $D(\{1, \dots, k^*\}) \in [\epsilon_1, 3\epsilon_1]$. As a result we may apply Claim 8, and we have that with probability at least $99/100$ there is an element $i_j \in \{i_1, \dots, i_\ell\}$ such that $i_j > k^*$ and i_j is ϵ_2 -heavy, i.e. $D(i_j) \geq (1 + \epsilon_2)D^*(i_j)$. We condition on this event going forward (the rest of our analysis will deal with this specific element i_j).

We now consider two cases:

Case 1: Distribution D has $D(\{1, \dots, i_j\}) \notin [1 - 3\epsilon_3, 1 + 3\epsilon_3]D^*(\{1, \dots, i_j\})$. Since the quantity $\widehat{D}(\{1, \dots, i_j\})$ obtained in Line 4 has expected value $D(\{1, \dots, i_j\}) \geq D(\{1, \dots, k^*\}) \geq \epsilon_1$, applying the multiplicative Chernoff bound implies that $\widehat{D}(\{1, \dots, i_j\}) \in [1 - \epsilon_3, 1 + \epsilon_3]D(\{1, \dots, i_j\})$ except with failure probability at most $\epsilon/10 \leq 1/10$. If this failure event does not occur then since $D(\{1, \dots, i_j\}) \notin [1 - 3\epsilon_3, 1 + 3\epsilon_3]D^*(\{1, \dots, i_j\})$ it must hold that $\widehat{D}(\{1, \dots, i_j\}) \notin [1 - \epsilon_3, 1 + \epsilon_3]D^*(\{1, \dots, i_j\})$ and consequently the algorithm outputs REJECT. Thus in Case 1 the algorithm outputs REJECT with overall failure probability at least $89/100$.

Case 2: Distribution D has $D(\{1, \dots, i_j\}) \in [1 - 3\epsilon_3, 1 + 3\epsilon_3]D^*(\{1, \dots, i_j\})$. This case is divided into two sub-cases depending on the value of $D^*(i_j)$.

Case 2(a): $D^*(i_j) \geq \epsilon_1$. In this case the algorithm reaches Line 6. We use the following claim:

Claim 12 *In Case 2(a), suppose that $i_j > k^*$ is such that $D(i_j) \geq (1 + \epsilon_2)D^*(i_j)$, and $D(\{1, \dots, i_j\}) \in [1 - 3\epsilon_3, 1 + 3\epsilon_3]D^*(\{1, \dots, i_j\})$. Then*

$$\frac{D(\{1, \dots, i_j - 1\})}{D(i_j)} \leq \left(1 - \frac{\epsilon_2}{4}\right) \cdot \frac{D^*(\{1, \dots, i_j - 1\})}{D^*(i_j)}.$$

Proof: To simplify notation we write

$$a \stackrel{\text{def}}{=} D(i_j); \quad b \stackrel{\text{def}}{=} D^*(i_j); \quad c \stackrel{\text{def}}{=} D(\{1, \dots, i_j - 1\}); \quad d \stackrel{\text{def}}{=} D^*(\{1, \dots, i_j - 1\}).$$

We have that

$$a \geq (1 + \epsilon_2)b \quad \text{and} \quad a + c \leq (1 + 3\epsilon_3)(b + d). \quad (44)$$

This gives

$$c \leq (1 + 3\epsilon_3)(b + d) - (1 + \epsilon_2)b = (1 + 3\epsilon_3)d + (3\epsilon_3 - \epsilon_2)b < (1 + 3\epsilon_3)d, \quad (45)$$

where in the last inequality we used $\epsilon_2 > 3\epsilon_3$. Recalling that $a \geq (1 + \epsilon_2)b$ and using $\epsilon_3 = \epsilon_2/24$ we get

$$\frac{c}{a} < \frac{(1 + 3\epsilon_3)d}{(1 + \epsilon_2)b} = \frac{d}{b} \cdot \frac{1 + \epsilon_2/8}{1 + \epsilon_2} < \frac{d}{b} \cdot \left(1 - \frac{\epsilon_2}{4}\right). \quad (46)$$

This proves the claim. \blacksquare

Applying Claim 12, we get that in Line 6 we have

$$\frac{D(\{1, \dots, i_j - 1\})}{D(i_j)} \leq \left(1 - \frac{\epsilon_2}{4}\right) \cdot \frac{D^*(\{1, \dots, i_j - 1\})}{D^*(i_j)}. \quad (47)$$

Recalling that by the premise of this case $D^*(i_j) \geq \epsilon_1$, by applying Claim 7 we have that $\{1, \dots, i_j - 1\}$ is an ϵ_1 -comparable witness for i_j . Therefore, by Lemma 2, with probability at least $1 - \frac{1}{10\ell}$ the call to COMPARE($\{i_j\}, \{1, \dots, i_j - 1\}$, $\frac{\epsilon_2}{16}, \frac{2}{\epsilon_1}, \frac{1}{10\ell}$) in Line 6 either outputs an element of $\{\text{High}, \text{Low}\}$ or outputs a value $v \leq (1 - \frac{\epsilon_2}{4})(1 + \frac{\epsilon_2}{16}) \frac{D^*(\{1, \dots, i_j - 1\})}{D^*(i_j)} < (1 - \frac{\epsilon_2}{8}) \frac{D^*(\{1, \dots, i_j - 1\})}{D^*(i_j)}$. In either case the algorithm outputs REJECT in Line 6, so we are done with Case 2(a).

Case 2(b): $D^*(i_j) < \epsilon_1$. In this case the algorithm reaches Line 10, and by item 2 of Claim 7, we have that S_1, \dots, S_t is a partition of $\{1, \dots, i_j - 1\}$ and each set S_1, \dots, S_t is a $\frac{1}{2}$ -comparable witness for i_j , i.e.,

$$\text{for all } i \in \{1, \dots, t\}, \quad \frac{1}{2}D^*(j) \leq D^*(S_i) \leq 2D^*(j). \quad (48)$$

We use the following lemma:

Claim 13 *In Case 2(b) suppose $i_j > k^*$ is such that $D(i_j) \geq (1 + \epsilon_2)D^*(i_j)$ and $D(\{1, \dots, i_j\}) \in [1 - 3\epsilon_3, 1 + 3\epsilon_3]D^*(\{1, \dots, i_j\})$. Then at least $(\epsilon_4/8)$ -fraction of the sets S_1, \dots, S_t are such that*

$$D(S_i) \leq (1 + \epsilon_4)D^*(S_i).$$

Proof: The proof is by contradiction. Let $\rho = 1 - \epsilon_4/8$ and suppose that there are w sets (without loss of generality we call them S_1, \dots, S_w) that satisfy $D(S_i) > (1 + \epsilon_4)D^*(S_i)$, where $\rho' = \frac{w}{t} > \rho$. We first observe that the weight of the w subsets S_1, \dots, S_w under D^* , as a fraction of $D^*({1, \dots, i_j - 1})$, is at least

$$\frac{D^*(S_1 \cup \dots \cup S_w)}{D^*(S_1 \cup \dots \cup S_w) + (t - w) \cdot 2D^*(j)} \geq \frac{w \frac{D^*(i_j)}{2}}{w \frac{D^*(i_j)}{2} + (t - w) \cdot 2D^*(j)} = \frac{w}{4t - 3w} = \frac{\rho'}{4 - 3\rho'},$$

where we used the right inequality in Equation (48) on S_{w+1}, \dots, S_t to obtain the leftmost expression above, and the left inequality in Equation (48) (together with the fact that $\frac{x}{x+c}$ is an increasing function of x for all $c > 0$) to obtain the inequality above. This implies that

$$\begin{aligned} D(\{1, \dots, i_j - 1\}) &= \sum_{i=1}^w D(S_i) + \sum_{i=w+1}^t D(S_i) \geq (1 + \epsilon_4) \sum_{i=1}^w D^*(S_i) + \sum_{i=w+1}^t D(S_i) \\ &\geq (1 + \epsilon_4) \frac{\rho'}{4 - 3\rho'} D^*({1, \dots, i_j - 1}) \\ &\geq (1 + \epsilon_4) \frac{\rho}{4 - 3\rho} D^*({1, \dots, i_j - 1}). \end{aligned} \quad (49)$$

From Equation (49) we have

$$\begin{aligned} D(\{1, \dots, i_j\}) &\geq (1 + \epsilon_4) \frac{\rho}{4 - 3\rho} D^*({1, \dots, i_j - 1}) + (1 + \epsilon_2) D^*(i_j) \\ &\geq \left(1 + \frac{3\epsilon_4}{8}\right) D^*({1, \dots, i_j - 1}) + (1 + \epsilon_2) D^*(i_j) \end{aligned}$$

where for the first inequality above we used $D(i_j) \geq (1 + \epsilon_2)D^*(i_j)$ and for the second inequality we used $(1 + \epsilon_4) \frac{\rho}{4 - 3\rho} \geq 1 + \frac{3\epsilon_4}{8}$. This implies that

$$D(\{1, \dots, i_j\}) \geq \left(1 + \frac{3\epsilon_4}{8}\right) D^*({1, \dots, i_j - 1}) + \left(1 + \frac{3\epsilon_4}{8}\right) D^*(i_j) = \left(1 + \frac{3\epsilon_4}{8}\right) D^*({1, \dots, i_j})$$

where the inequality follows from $\epsilon_2 \geq \frac{3\epsilon_4}{8}$. Since $\frac{3\epsilon_4}{8} > 3\epsilon_3$, though, this is a contradiction and the claim is proved. ■

Applying Claim 13, and recalling that $h = \Theta(1/\epsilon) = \Theta(1/\epsilon_4)$ sets are chosen randomly in Line 9, we have that with probability at least 9/10 there is some $r \in \{1, \dots, h\}$ such that $D(S_{a_r}) \leq (1 + \epsilon_4)D^*(S_{a_r})$. Combining this with $D(i_j) \geq (1 + \epsilon_2)D^*(i_j)$, we get that

$$\frac{D(S_{a_r})}{D(i_j)} \leq \frac{1 + \epsilon_4}{1 + \epsilon_2} \cdot \frac{D^*(S_{a_r})}{D^*(i_j)} \leq \left(1 - \frac{\epsilon_4}{2}\right) \cdot \frac{D^*(S_{a_r})}{D^*(i_j)}.$$

By Lemma 2, with probability at least $1 - \frac{1}{10\ell h}$ the call to COMPARE($\{i_j\}, S_{a_r}, \frac{\epsilon_4}{8}, 4, \frac{1}{10\ell n}$) in Line 10 either outputs an element of $\{\text{High}, \text{Low}\}$ or outputs a value $v \leq (1 - \frac{\epsilon_4}{2})(1 + \frac{\epsilon_4}{8}) \frac{D^*(S_{a_r})}{D^*(i_j)} < (1 - \frac{\epsilon_4}{4}) \frac{D^*(S_{a_r})}{D^*(i_j)}$. In either case the algorithm outputs REJECT in Line 10, so we are done in Case 2(b). This concludes the proof of soundness and the proof of Theorem 6. ■

6 Testing equality between two unknown distributions

6.1 An approach based on PCOND queries

In this subsection we consider the problem of testing whether two unknown distributions D_1, D_2 are identical versus ϵ -far, given PCOND access to these distributions. Although this is known to require $\Omega(N^{2/3})$ many samples in the standard model [BFR⁺10, Val11], we are able to give a $\text{poly}(\log N, 1/\epsilon)$ -query algorithm using PCOND queries, by taking advantage of comparisons to perform some sort of *clustering* of the domain.

On a high level the algorithm works as follows. First it obtains (with high probability) a small set of points R such that almost every element in $[N]$, except possibly for some negligible subset according to D_1 , has probability weight (under D_1) close to some “representative” in R . Next, for each representative r in R it obtains an estimate of the weight, according to D_1 , of a set of points U such that $D_1(u)$ is close to $D_1(r)$ for each u in U (i.e. r ’s “neighborhood under D_1 ”). This is done using the procedure ESTIMATE-NEIGHBORHOOD from Subsection 3.2). Note that these neighborhoods can be interpreted roughly as a succinct *cover* of the support of D_1 into (not necessarily disjoint) sets of points, where within each set the points have similar weight (according to D_1). Our algorithm is based on the observation that, if D_1 and D_2 are far from each other, it must be the case that one of these sets, denoted U^* , reflects it in one of the following ways: (1) $D_2(U^*)$ differs significantly from $D_1(U^*)$; (2) U^* contains a subset of points V^* such that $D_2(v)$ differs significantly from $D_2(r)$ for each v in V^* , and either $D_1(V^*)$ is relatively large or $D_2(V^*)$ is relatively large. (This structural result is made precise in Lemma 15). We thus take additional samples, both from D_1 and from D_2 , and compare the weight (according to both distributions) of each point in these samples to the representatives in R (using the procedure COMPARE from Subsection 3.1). In this manner we detect (with high probability) that either (1) or (2) holds.

We begin by formalizing the notion of a cover discussed above:

Definition 5 (Weight-Cover) *Given a distribution D on $[N]$ and a parameter $\epsilon_1 > 0$, we say that a point $i \in [N]$ is ϵ_1 -covered by a set $R = \{r_1, \dots, r_t\} \subseteq [N]$ if there exists a point $r_j \in R$ such that $D(i) \in [1/(1 + \epsilon_1), 1 + \epsilon_1]D(r_j)$. Let the set of points in $[N]$ that are ϵ_1 -covered by R be denoted by $C_{\epsilon_1}^D(R)$. We say that R is an (ϵ_1, ϵ_2) -cover for D if $D([N] \setminus C_{\epsilon_1}^D(R)) \leq \epsilon_2$.*

The following lemma says that a small sample of points drawn from D gives a cover with high probability:

Lemma 14 *Let D be any distribution over $[N]$. Given any fixed $c > 0$, there exists a constant $c' > 0$ such that with probability at least 99/100, a sample R of size $m = c' \frac{\log(N/\epsilon)}{\epsilon^2} \cdot \log\left(\frac{\log(N/\epsilon)}{\epsilon}\right)$ drawn according to distribution D is an $(\epsilon/c, \epsilon/c)$ -cover for D .*

Proof: Let t denote $\lceil \ln(2cN/\epsilon) \cdot \frac{c}{\epsilon} \rceil$. We define t “buckets” of points with similar weight under D as follows: for $i = 0, 1, \dots, t - 1$, define $B_i \subseteq [N]$ to be

$$B_i \stackrel{\text{def}}{=} \left\{ x \in [N] : \frac{1}{(1 + \epsilon/c)^{i+1}} < D(x) \leq \frac{1}{(1 + \epsilon/c)^i} \right\}.$$

Let L be the set of points x which are not in any of B_0, \dots, B_{t-1} (because $D(x)$ is too small); since every point in L has $D(x) < \frac{\epsilon}{2cN}$, one can see that $D(L) \leq \frac{\epsilon}{2c}$.

It is easy to see that if the sample R contains a point from a bucket B_j then every point $y \in B_j$ is $\frac{\epsilon}{c}$ -covered by R . We say that bucket B_i is *insignificant* if $D(B_i) \leq \frac{\epsilon}{2ct}$; otherwise bucket B_i is *significant*. It is clear that the total weight under D of all insignificant buckets is at most $\epsilon/2c$. Thus if we can show that for the claimed sample size, with probability at least 99/100 every significant bucket has at least one of its points in R , we will have established the lemma.

This is a simple probabilistic calculation: fix any significant bucket B_j . The probability that m random draws from D all miss B_j is at most $(1 - \frac{\epsilon}{2ct})^m$, which is at most $\frac{1}{100t}$ for a suitable (absolute constant) choice of c' . Thus a union bound over all (at most t) significant buckets gives that with probability at least 99/100, no significant bucket is missed by R . ■

Lemma 15 *Suppose $d_{\text{TV}}(D_1, D_2) \geq \epsilon$, and let $R = \{r_1, \dots, r_t\}$ be an $(\tilde{\epsilon}, \tilde{\epsilon})$ -cover for D_1 where $\tilde{\epsilon} \leq \epsilon/100$. Then, there exists $j \in [t]$ such that at least one of the following conditions holds for every $\alpha \in [\tilde{\epsilon}, 2\tilde{\epsilon}]$:*

1. $D_1(U_\alpha^{D_1}(r_j)) \geq \frac{\tilde{\epsilon}}{t}$ and $D_2(U_\alpha^{D_1}(r_j)) \notin [1 - \tilde{\epsilon}, 1 + \tilde{\epsilon}]D_1(U_\alpha^{D_1}(r_j))$, or $D_1(U_\alpha^{D_1}(r_j)) < \frac{\tilde{\epsilon}}{t}$ and $D_2(U_\alpha^{D_1}(r_j)) > \frac{2\tilde{\epsilon}}{t}$;
2. $D_1(U_\alpha^{D_1}(r_j)) \geq \frac{\tilde{\epsilon}}{t}$, and at least a $\tilde{\epsilon}$ -fraction of the points i in $U_\alpha^{D_1}(r_j)$ satisfy $\frac{D_2(i)}{D_2(r_j)} \notin [1/(1 + \alpha + \tilde{\epsilon}), 1 + \alpha + \tilde{\epsilon}]$;
3. $D_1(U_\alpha^{D_1}(r_j)) \geq \frac{\tilde{\epsilon}}{t}$, and the total weight according to D_2 of the points i in $U_\alpha^{D_1}(r_j)$ for which $\frac{D_2(i)}{D_2(r_j)} \notin [1/(1 + \alpha + \tilde{\epsilon}), 1 + \alpha + \tilde{\epsilon}]$ is at least $\frac{\tilde{\epsilon}^2}{t}$;

Proof: Without loss of generality, we can assume that $\epsilon \leq 1/4$. Suppose, contrary to the claim, that for each r_j there exists $\alpha_j \in [\tilde{\epsilon}, 2\tilde{\epsilon}]$ such that if we let $U_j \stackrel{\text{def}}{=} U_{\alpha_j}^{D_1}(r_j)$, then the following holds:

1. If $D_1(U_j) < \frac{\tilde{\epsilon}}{t}$, then $D_2(U_j) \leq \frac{2\tilde{\epsilon}}{t}$;
2. If $D_1(U_j) \geq \frac{\tilde{\epsilon}}{t}$, then:
 - (a) $D_2(U_j) \in [1 - \tilde{\epsilon}, 1 + \tilde{\epsilon}]D_1(U_j)$;
 - (b) Less than an $\tilde{\epsilon}$ -fraction of the points y in U_j satisfy $\frac{D_2(y)}{D_2(r_j)} \notin [1/(1 + \alpha_j + \tilde{\epsilon}), 1 + \alpha_j + \tilde{\epsilon}]$;
 - (c) The total weight according to D_2 of the points y in U_j for which $\frac{D_2(y)}{D_2(r_j)} \notin [1/(1 + \alpha_j + \tilde{\epsilon}), 1 + \alpha_j + \tilde{\epsilon}]$ is at most $\frac{\tilde{\epsilon}^2}{t}$;

We show that in such a case $d_{\text{TV}}(D_1, D_2) < \epsilon$, contrary to the premise of the claim.

Consider each point $r_j \in R$ such that $D_1(U_j) \geq \frac{\tilde{\epsilon}}{t}$. By the foregoing discussion (point 2(a)), $D_2(U_j) \in [1 - \tilde{\epsilon}, 1 + \tilde{\epsilon}]D_1(U_j)$. By the definition of U_j (and since $\alpha_j \leq 2\tilde{\epsilon}$),

$$D_1(r_j) \in [1/(1 + 2\tilde{\epsilon}), 1 + 2\tilde{\epsilon}] \frac{D_1(U_j)}{|U_j|}. \quad (50)$$

Turning to bound $D_2(r_j)$, on one hand (by 2(b))

$$D_2(U_j) = \sum_{y \in U_j} D_2(y) \geq \tilde{\epsilon}|U_j| \cdot 0 + (1 - \tilde{\epsilon})|U_j| \cdot \frac{D_2(r_j)}{1 + 3\tilde{\epsilon}}, \quad (51)$$

and so

$$D_2(r_j) \leq \frac{(1 + 3\tilde{\epsilon})D_2(U_j)}{(1 - \tilde{\epsilon})|U_j|} \leq (1 + 6\tilde{\epsilon}) \frac{D_1(U_j)}{|U_j|}. \quad (52)$$

On the other hand (by 2(c)),

$$D_2(U_j) = \sum_{y \in U_j} D_2(y) \leq \frac{\tilde{\epsilon}^2}{t} + |U_j| \cdot (1 + 3\tilde{\epsilon})D_2(r_j), \quad (53)$$

and so

$$D_2(r_j) \geq \frac{D_2(U_j) - \tilde{\epsilon}^2/t}{(1 + 3\tilde{\epsilon})|U_j|} \geq \frac{(1 - \tilde{\epsilon})D_1(U_j) - \tilde{\epsilon}D_1(U_j)}{(1 + 3\tilde{\epsilon})|U_j|} \geq (1 - 5\tilde{\epsilon}) \frac{D_1(U_j)}{|U_j|}. \quad (54)$$

Therefore, for each such r_j we have

$$D_2(r_j) \in [1 - 8\tilde{\epsilon}, 1 + 10\tilde{\epsilon}]D_1(r_j). \quad (55)$$

Let $C \stackrel{\text{def}}{=} \bigcup_{j=1}^t U_j$. We next partition the points in C so that each point $i \in C$ is assigned to some $r_{j(i)}$ such that $i \in U_{j(i)}$. We define the following ‘‘bad’’ subsets of points in $[N]$:

1. $B_1 \stackrel{\text{def}}{=} [N] \setminus C$, so that $D_1(B_1) \leq \tilde{\epsilon}$ (we later bound $D_2(B_1)$);
2. $B_2 \stackrel{\text{def}}{=} \{i \in C : D_1(U_{j(i)}) < \tilde{\epsilon}/t\}$, so that $D_1(B_2) \leq \tilde{\epsilon}$ and $D_2(B_2) \leq 2\tilde{\epsilon}$;
3. $B_3 \stackrel{\text{def}}{=} \{i \in C \setminus B_2 : D_2(i) \notin [1/(1 + 3\tilde{\epsilon}), 1 + 3\tilde{\epsilon}]D_2(r_{j(i)})\}$, so that $D_1(B_3) \leq 2\tilde{\epsilon}$ and $D_2(B_3) \leq \tilde{\epsilon}^2$.

Let $B \stackrel{\text{def}}{=} B_1 \cup B_2 \cup B_3$. Observe that for each $i \in [N] \setminus B$ we have that

$$D_2(i) \in [1/(1 + 3\tilde{\epsilon}), 1 + 3\tilde{\epsilon}]D_2(r_{j(i)}) \subset [1 - 15\tilde{\epsilon}, 1 + 15\tilde{\epsilon}]D_1(r_{j(i)}) \subset [1 - 23\tilde{\epsilon}, 1 + 23\tilde{\epsilon}]D_1(i), \quad (56)$$

where the first containment follows from the fact that $i \notin B$, the second follows from Equation (55), and the third from the fact that $i \in U_{j(i)}$. In order to complete the proof we need a bound on $D_2(B_1)$, which we obtain next.

$$\begin{aligned} D_2(B_1) &= 1 - D_2([N] \setminus B_1) \leq 1 - D_2([N] \setminus B) \leq 1 - (1 - 23\tilde{\epsilon})D_1([N] \setminus B) \\ &\leq 1 - (1 - 23\tilde{\epsilon})(1 - 4\tilde{\epsilon}) \leq 27\tilde{\epsilon}. \end{aligned} \quad (57)$$

Therefore,

$$\begin{aligned} d_{\text{TV}}(D_1, D_2) &= \frac{1}{2} \sum_{i=1}^N |D_1(i) - D_2(i)| \\ &\leq \frac{1}{2} \left(D_1(B) + D_2(B) + \sum_{i \notin B} 23\tilde{\epsilon}D_1(i) \right) < \epsilon, \end{aligned} \quad (58)$$

and we have reached a contradiction. \blacksquare

Algorithm 8: Algorithm PCOND $_{D_1, D_2}$ -TEST-EQUALITY-UNKNOWN

Input: PCOND query access to distributions D_1 and D_2 and a parameter ϵ .

1. Set $\tilde{\epsilon} = \epsilon/100$.
 2. Draw a sample R of size $t = \tilde{\Theta}\left(\frac{\log N}{\epsilon^2}\right)$ from D_1 .
 3. For each $r_j \in R$:
 - (a) Call ESTIMATE-NEIGHBORHOOD $_{D_1}$ on r_j with $\kappa = \tilde{\epsilon}$, $\eta = \frac{\tilde{\epsilon}}{8}$, $\beta = \frac{\tilde{\epsilon}}{2t}$, $\delta = \frac{1}{100t}$ and let the output be denoted by $(\hat{w}_j^{(1)}, \alpha_j)$.
 - (b) Set $\theta = \kappa\eta\beta\delta/64 = \tilde{\Theta}(\epsilon^7/\log^2 N)$.
 - (c) Draw a sample S_1 from D_1 , of size $s_1 = \Theta\left(\frac{t}{\epsilon^2}\right) = \tilde{\Theta}\left(\frac{\log N}{\epsilon^4}\right)$.
 - (d) Draw a sample S_2 from D_2 , of size $s_2 = \Theta\left(\frac{t \log t}{\epsilon^3}\right) = \tilde{\Theta}\left(\frac{\log N}{\epsilon^5}\right)$.
 - (e) For each point $i \in S_1 \cup S_2$ call COMPARE $_{D_1}(\{r_j\}, \{i\}, \theta/4, 4, 1/(200t(s_1 + s_2)))$ and COMPARE $_{D_2}(\{r_j\}, \{i\}, \theta/4, 4, 1/(200t(s_1 + s_2)))$, and let the outputs be denoted $\rho_{r_j}^{(1)}(i)$ and $\rho_{r_j}^{(2)}(i)$, respectively (where in particular these outputs may be High or Low).
 - (f) Let $\hat{w}_j^{(2)}$ be the fraction of occurrences of $i \in S_2$ such that $\rho_{r_j}^{(1)}(i) \in [1/(1 + \alpha_j + \theta/2), 1 + \alpha_j + \theta/2]$.
 - (g) If $(\hat{w}_j^{(1)} \leq \frac{3}{4}\tilde{\epsilon}$ and $\hat{w}_j^{(2)} > \frac{3}{2}\tilde{\epsilon})$ or $(\hat{w}_j^{(1)} > \frac{3}{4}\tilde{\epsilon}$ and $\hat{w}_j^{(2)}/\hat{w}_j^{(1)} \notin [1 - \tilde{\epsilon}/2, 1 + \tilde{\epsilon}/2])$, then output REJECT.
 - (h) If there exists $i \in S_1 \cup S_2$ such that $\rho_{r_j}^{(1)}(i) \in [1/(\alpha_j + \tilde{\epsilon}/2), 1 + \alpha_j + \tilde{\epsilon}/2]$ and $\rho_{r_j}^{(2)}(i) \notin [1/(\alpha_j + 3\tilde{\epsilon}/2), 1 + \alpha_j + 3\tilde{\epsilon}/2]$, then output REJECT.
 4. Output ACCEPT.
-

Theorem 10 *If $D_1 = D_2$ then with probability at least $2/3$ Algorithm PCOND-TEST-EQUALITY-UNKNOWN returns ACCEPT, and if $d_{TV}(D_1, D_2) \geq \epsilon$, then with probability at least $2/3$ Algorithm PCOND-TEST-EQUALITY-UNKNOWN returns REJECT. The number of PCOND queries performed by the algorithm is $\tilde{O}\left(\frac{\log^6 N}{\epsilon^{21}}\right)$.*

Proof: The number of queries performed by the algorithm is the sum of: (1) t times the number of queries performed in each execution of ESTIMATE-NEIGHBORHOOD (in Line 3-a) and (2) $t \cdot (s_1 + s_2) = O(t \cdot s_2)$ times the number of queries performed in each execution of COMPARE (in Line 3-e). By Lemma 3 (and the settings of the parameters in the calls to ESTIMATE-NEIGHBORHOOD), the first term is $O\left(t \cdot \frac{\log(1/\delta) \cdot \log(\log(1/\delta)/(\beta\eta^2))}{\kappa^2\eta^4\beta^3\delta^2}\right) = \tilde{O}\left(\frac{\log^6 N}{\epsilon^{19}}\right)$, and by Lemma 2 (and the settings of the parameters in the calls to COMPARE), the second term is $O\left(t \cdot s_2 \cdot \frac{\log(t \cdot s_2)}{\theta^2}\right) = \tilde{O}\left(\frac{\log^6 N}{\epsilon^{21}}\right)$, so that we get the bound stated in the theorem.

We now turn to establishing the correctness of the algorithm. We shall use the shorthand U_j for

$U_{\alpha_j}^{D_1}(r_j)$, and U'_j for $U_{\alpha_j+\theta}^{D_1}(r_j)$. We consider the following “desirable” events.

1. The event E_1 is that the sample R is a $(\tilde{\epsilon}, \tilde{\epsilon})$ -weight-cover for D_1 (for $\tilde{\epsilon} = \epsilon/100$). By Lemma 14 (and an appropriate constant in the $\Theta(\cdot)$ notation for the size of R), the probability that E_1 holds is at least $99/100$.
2. The event E_2 is that all calls to the procedure ESTIMATE-NEIGHBORHOOD are as specified by Lemma 3. By the setting of the confidence parameter in the calls to the procedure, the event E_2 holds with probability at least $99/100$.
3. The event E_3 is that all calls to the procedure COMPARE are as specified by Lemma 2. By the setting of the confidence parameter in the calls to the procedure, the event E_3 holds with probability at least $99/100$.
4. The event E_4 is that $D_2(U'_j \setminus U_j) \leq \eta\beta/16 = \tilde{\epsilon}^2/(256t)$ for each j . If $D_2 = D_1$ then this event follows from E_2 . Otherwise, it holds with probability at least $99/100$ by the setting of θ and the choice of α_j (as shown in the proof of Lemma 3 in the analysis of the event E_1 there).
5. The event E_5 is defined as follows. For each j , if $D_2(U_j) \geq \tilde{\epsilon}/(4t)$, then $|S_2 \cap U_j|/|S_2| \in [1 - \tilde{\epsilon}/10, 1 + \tilde{\epsilon}/10]D_2(U_j)$, and if $D_2(U_j) < \tilde{\epsilon}/(4t)$ then $|S_2 \cap U_j|/|S_2| < (1 + \tilde{\epsilon}/10)\tilde{\epsilon}/(4t)$. This event holds with probability at least $99/100$ by applying a multiplicative Chernoff bound in the first case, and Corollary 2 in the second.
6. The event E_6 is that for each j we have $|S_2 \cap (U'_j \setminus U_j)|/|S_2| \leq \tilde{\epsilon}^2/(128t)$. Conditioned on E_4 , the event E_6 holds with probability at least $99/100$ by applying Corollary 2.

From this point on we assume that events $E_1 - E_6$ all hold. Note that in particular this implies the following:

1. By E_2 , for every j :
 - If $D_1(U_j) \geq \beta = \tilde{\epsilon}/(2t)$, then $\hat{w}_j^{(1)} \in [1 - \eta, 1 + \eta]D_1(U_j) = [1 - \tilde{\epsilon}/8, 1 + \tilde{\epsilon}/8]D_1(U_j)$.
 - If $D_1(U_j) < \tilde{\epsilon}/(2t)$, then $\hat{w}_j^{(1)} \leq (1 + \tilde{\epsilon}/8)(\tilde{\epsilon}/(2t))$.
2. By E_3 , for every j and for each point $i \in S_1 \cup S_2$:
 - If $i \in U_j$, then $\rho_{r_j}^{(1)}(i) \in [1/(1 + \alpha_j + \frac{\theta}{2}), 1 + \alpha_j + \frac{\theta}{2}]$.
 - If $i \notin U'_j$, then $\rho_{r_j}^{(1)}(i) \notin [1/(1 + \alpha_j + \frac{\theta}{2}), 1 + \alpha_j + \frac{\theta}{2}]$.
3. By the previous item and $E_4 - E_6$:
 - If $D_2(U_j) \geq \tilde{\epsilon}/(4t)$, then $\hat{w}_j^{(2)} \geq (1 - \tilde{\epsilon}/10)D_2(U_j)$ and $\hat{w}_j^{(2)} \leq (1 + \tilde{\epsilon}/10)D_2(U_j) + \tilde{\epsilon}^2/(128t) \leq (1 + \tilde{\epsilon}/8)D_2(U_j)$.
 - If $D_2(U_j) < \tilde{\epsilon}/(4t)$ then $\hat{w}_j^{(2)} \leq (1 + \tilde{\epsilon}/10)\tilde{\epsilon}/(4t) + \tilde{\epsilon}^2/(128t) \leq (1 + \tilde{\epsilon}/4)(\tilde{\epsilon}/(4t))$.

Completeness. Assume D_1 and D_2 are the same distribution D . For each j , if $D(U_j) \geq \tilde{\epsilon}/t$, then by the foregoing discussion, $\hat{w}_j^{(1)} \geq (1 - \tilde{\epsilon}/8)D(U_j) > 3\tilde{\epsilon}/(4t)$ and $\hat{w}_j^{(2)}/\hat{w}_j^{(1)} \in [(1 - \tilde{\epsilon}/8)^2, (1 + \tilde{\epsilon}/8)^2] \subset [1 - \tilde{\epsilon}/2, 1 + \tilde{\epsilon}/2]$, so that the algorithm does not reject in Line 3-g. Otherwise (i.e., $D(U_j) < \tilde{\epsilon}/t$), we consider two subcases. Either $D(U_j) \leq \tilde{\epsilon}/(2t)$, in which case $\hat{w}_j^{(1)} \leq 3\tilde{\epsilon}/(4t)$, or $\tilde{\epsilon}/(2t) < D(U_j) < \tilde{\epsilon}/t$, and then $\hat{w}_j^{(1)} \in [1 - \tilde{\epsilon}/8, 1 + \tilde{\epsilon}/8]D_1(U_j)$. Since in both cases $\hat{w}_j^{(2)} \leq (1 + \tilde{\epsilon}/8)D(U_j) \leq 3\tilde{\epsilon}/(2t)$, the algorithm does not reject in Line 3-g. By E_3 , the algorithm does not reject in Line 3-h either. We next turn to establish soundness.

Soundness. Assume $d_{\text{TV}}(D_1, D_2) \geq \epsilon$. By applying Lemma 15 on R (and using E_1), there exists an index j for which one of the items in the lemma holds. We denote this index by j^* , and consider the three items in the lemma.

1. If Item 1 holds, then we consider its two cases:

- (a) In the first case, $D_1(U_{j^*}) \geq \tilde{\epsilon}/t$ and $D_2(U_{j^*}) \notin [1 - \tilde{\epsilon}, 1 + \tilde{\epsilon}]D_1(U_{j^*})$. Due to the lower bound on $D_1(U_{j^*})$ we have that $\hat{w}_{j^*}^{(1)} \in [1 - \tilde{\epsilon}/8, 1 + \tilde{\epsilon}/8]D_1(U_{j^*})$, so that in particular $\hat{w}_{j^*}^{(1)} > 3\tilde{\epsilon}/(4t)$. As for $\hat{w}_{j^*}^{(2)}$, either $\hat{w}_{j^*}^{(2)} < (1 - \tilde{\epsilon})(1 + \tilde{\epsilon}/8)D_1(U_{j^*})$ (this holds both when $D_2(U_{j^*}) \geq \tilde{\epsilon}/(4t)$ and when $D_2(U_{j^*}) < \tilde{\epsilon}/(4t)$) or $\hat{w}_{j^*}^{(2)} > (1 + \tilde{\epsilon})(1 - \tilde{\epsilon}/10)D_1(U_{j^*})$. In either (sub)case $\hat{w}_{j^*}^{(2)}/\hat{w}_{j^*}^{(1)} \notin [1 - \tilde{\epsilon}/2, 1 + \tilde{\epsilon}/2]$, causing the algorithm to reject in (the second part of) Line 3-g.
- (b) In the second case, $D_1(U_{j^*}) < \tilde{\epsilon}/t$ and $D_2(U_{j^*}) > 2\tilde{\epsilon}/t$. Due to the lower bound on $D_2(U_{j^*})$ we have that $\hat{w}_{j^*}^{(2)} \geq (1 - \tilde{\epsilon}/10)D_2(U_{j^*}) > (1 - \tilde{\epsilon}/10)(2\tilde{\epsilon}/t)$, so that in particular $\hat{w}_{j^*}^{(2)} > (3\tilde{\epsilon}/(2t))$. As for $\hat{w}_{j^*}^{(1)}$, if $D_1(U_{j^*}) \leq \tilde{\epsilon}/(2t)$, then $\hat{w}_{j^*}^{(1)} \leq 3\tilde{\epsilon}/(4t)$, causing the algorithm to reject in (the first part of) Line 3-g. If $\tilde{\epsilon}/(2t) < D_1(U_{j^*}) \leq \tilde{\epsilon}/t$, then $\hat{w}_{j^*}^{(1)} \in [1 - \tilde{\epsilon}/8, 1 + \tilde{\epsilon}/8]D_1(U_{j^*}) \leq (1 + \tilde{\epsilon}/8)(\tilde{\epsilon}/t)$, so that $\hat{w}_{j^*}^{(2)}/\hat{w}_{j^*}^{(1)} \geq \frac{(1 - \tilde{\epsilon}/10)(2\tilde{\epsilon}/t)}{(1 + \tilde{\epsilon}/8)\tilde{\epsilon}/t} > (1 + \tilde{\epsilon}/2)$, causing the algorithm to reject in (either the first or second part of) Line 3-g.

2. If Item 2 holds, then by the choice of the size of S_1 , which is $\Theta(t/\tilde{\epsilon}^2)$, with probability at least 99/100, the sample S_1 will contain a point i for which $\frac{D_2(i)}{D_2(r_{j^*})} \notin [1/(1 + \alpha_{j^*} + \tilde{\epsilon}), 1 + \alpha_{j^*} + \tilde{\epsilon}]$, and by E_3 this will be detected in Line 3-h.
3. Similarly, if Item 3 holds, then by the choice of the size of S_2 , with probability at least 99/100, the sample S_2 will contain a point i for which $\frac{D_2(i)}{D_2(r_{j^*})} \notin [1/(1 + \alpha_{j^*} + \tilde{\epsilon}), 1 + \alpha_{j^*} + \tilde{\epsilon}]$, and by E_3 this will be detected in Line 3-h.

The theorem is thus established. ■

6.2 An approach based on simulating EVAL

In this subsection we present an alternate approach for testing whether two unknown distributions D_1, D_2 are identical versus ϵ -far. We prove the following theorem:

Theorem 11 COND-TEST-EQUALITY-UNKNOWN *is a*

$$\tilde{O}\left(\frac{(\log N)^5 \cdot (\log(1/\delta))^2}{\epsilon^4}\right)$$

-query algorithm with the following properties: given $\text{COND}_{D_1}, \text{COND}_{D_2}$ oracles for any two distributions D_1, D_2 over $[N]$, it outputs **ACCEPT** with probability at least $2/3$ if $D_1 = D_2$ and outputs **REJECT** with probability at least $2/3$ if $d_{\text{TV}}(D_1, D_2) \geq \epsilon$.

At the heart of this result is an efficient simulation of an “approximate EVAL_D oracle” using a COND_D oracle. (Recall that an EVAL_D oracle is an oracle which, given as input an element $i \in [N]$, outputs the numerical value $D(i)$.) We feel that this efficient simulation of an approximate EVAL oracle using a COND oracle is of independent interest since it sheds light on the relative power of the COND and EVAL models.

In more detail, the starting point of our approach to prove Theorem 11 is a simple algorithm from [RS09] that uses an EVAL_D oracle to test equality between D and a known distribution D^* . We first show (see Theorem 12) that a modified version of the algorithm, which uses a **SAMP** oracle and an “approximate” EVAL oracle, can be used to efficiently test equality between two unknown distributions D_1 and D_2 . We then show (see Theorem 13) how the required “approximate” EVAL oracle can be efficiently implemented using a COND oracle. Theorem 11 follows straightforwardly by combining Theorems 12 and 13.

6.2.1 Approximate EVAL oracles.

We begin by defining the notion of an “approximate EVAL oracle” that we will use. Intuitively this is an oracle which gives a multiplicatively $(1 \pm \epsilon)$ -accurate estimate of the value of $D(i)$ for all i in a fixed set of probability weight at least $1 - \epsilon$ under D . More precisely, we have the following definition:

Definition 6 Let D be a distribution over $[N]$. An (ϵ, δ) -approximate EVAL_D simulator is a randomized procedure **ORACLE** with the following property: For each $0 < \epsilon < 1$, there is a fixed set $S^{(\epsilon, D)} \subsetneq [N]$ with $D(S^{(\epsilon, D)}) < \epsilon$ for which the following holds. Given as input an element $i^* \in [N]$, the procedure **ORACLE** either outputs a value $\alpha \in [0, 1]$ or outputs **UNKNOWN**. The following holds for all $i^* \in [N]$:

- (i) If $i^* \notin S^{(\epsilon, D)}$ then with probability at least $1 - \delta$ the output of **ORACLE** on input i^* is a value $\alpha \in [0, 1]$ such that $\alpha \in [1 - \epsilon, 1 + \epsilon]D(i^*)$;
- (i) If $i^* \in S^{(\epsilon, D)}$ then with probability at least $1 - \delta$ the procedure either outputs **UNKNOWN** or outputs a value $\alpha \in [0, 1]$ such that $\alpha \in [1 - \epsilon, 1 + \epsilon]D(i^*)$.

We note that according to the above definition, it may be the case that different calls to **ORACLE** on the same input element $i^* \in [N]$ may return different values. However, the “low-weight” set $S^{(\epsilon, D)}$ is an *a priori* fixed set that does not depend in any way on the input point i^* given to the algorithm. The key property of an (ϵ, δ) -approximate EVAL_D oracle is that it reliably gives a multiplicatively $(1 \pm \epsilon)$ -accurate estimate of the value of $D(i)$ for all i in some fixed set of probability weight at least $1 - \epsilon$ under D .

6.2.2 Testing equality between D_1 and D_2 using an approximate EVAL oracle.

We now show how an approximate EVAL_{D_1} oracle, an approximate EVAL_{D_2} oracle, and a SAMP_{D_1} oracle can be used together to test whether $D_1 = D_2$ versus $d_{\text{TV}}(D_1, D_2) \geq \epsilon$. As mentioned earlier, the approach is a simple extension of the EVAL algorithm given in Observation 24 of [RS09].

Theorem 12 *Let ORACLE_1 be an $(\epsilon/100, \epsilon/100)$ -approximate EVAL_{D_1} simulator and let ORACLE_2 be an $(\epsilon/100, \epsilon/100)$ -approximate EVAL_{D_2} simulator. There is an algorithm $\text{TEST-EQUALITY-UNKNOWN}$ with the following properties: for any distributions D_1, D_2 over $[N]$, algorithm $\text{TEST-EQUALITY-UNKNOWN}$ makes $O(1/\epsilon)$ queries to ORACLE_1 , ORACLE_2 and SAMP_{D_1} , and it outputs ACCEPT with probability at least $7/10$ if $D_1 = D_2$ and outputs REJECT with probability at least $7/10$ if $d_{\text{TV}}(D_1, D_2) \geq \epsilon$.*

Algorithm 9: TEST-EQUALITY-UNKNOWN

Input: query access to ORACLE_1 , to ORACLE_2 , and access to SAMP_{D_1} oracle

- 1: Call the SAMP_{D_1} oracle $m = 5/\epsilon$ times to obtain points h_1, \dots, h_m distributed according to D_1 .
- 2: Call the SAMP_{D_2} oracle $m = 5/\epsilon$ times to obtain points h_{m+1}, \dots, h_{2m} distributed according to D_2 .
- 3: **for** $j = 1$ to $2m$ **do**
- 4: Call $\text{ORACLE}_1(h_j)$. If it returns UNKNOWN then output REJECT , otherwise let $v_{1,i} \in [0, 1]$ be the value it outputs.
- 5: Call $\text{ORACLE}_2(h_j)$. If it returns UNKNOWN then output REJECT , otherwise let $v_{2,i} \in [0, 1]$ be the value it outputs.
- 6: **if** $v_{1,j} \notin [1 - \epsilon/8, 1 + \epsilon/8]v_{2,j}$ **then**
- 7: output REJECT and exit
- 8: **end if**
- 9: **end for**
- 10: output ACCEPT

It is clear that $\text{TEST-EQUALITY-UNKNOWN}$ makes $O(1/\epsilon)$ queries as claimed. To prove Theorem 12 we argue completeness and soundness below.

Completeness: Suppose that $D_1 = D_2$. Since ORACLE_1 is an $(\epsilon/100, \epsilon/100)$ -approximate EVAL_{D_1} simulator, the probability that any of the $2m = 10/\epsilon$ points h_1, \dots, h_{2m} drawn in Lines 1 and 2 lies in $S^{(\epsilon/100, D_1)}$ is at most $1/10$. Going forth, let us assume that all points h_i indeed lie outside $S^{(\epsilon/100, D_1)}$. Then for each execution of Line 4 we have that with probability at least $1 - \epsilon/100$ the call to $\text{ORACLE}(h_i)$ yields a value $v_{1,i}$ satisfying $v_{1,i} \in [1 - \frac{\epsilon}{100}, 1 + \frac{\epsilon}{100}]D_1(i)$. The same holds for each execution of Line 5. Since there are $20/\epsilon$ total executions of Lines 4 and 5, with overall probability at least $7/10$ we have that each $1 \leq j \leq m$ has $v_{1,j}, v_{2,j} \in [1 - \frac{\epsilon}{100}, 1 + \frac{\epsilon}{100}]D_1(i)$. If this is the case then $v_{1,j}, v_{2,j}$ pass the check in Line 6, and thus the algorithm outputs ACCEPT with overall probability at least $7/10$.

Soundness: Now suppose that $d_{\text{TV}}(D_1, D_2) \geq \epsilon$. Let us say that $i \in [N]$ is *good* if $D_1(i) \in [1 - \epsilon/5, 1 + \epsilon/5]D_2(i)$. Let $\text{BAD} \subseteq [N]$ denote the set of all $i \in [N]$ that are not good. We have

$$2d_{\text{TV}}(D_1, D_2) = \sum_{i \text{ is good}} |D_1(i) - D_2(i)| + \sum_{i \text{ is bad}} |D_1(i) - D_2(i)| \geq 2\epsilon.$$

Since

$$\sum_{i \text{ is good}} |D_1(i) - D_2(i)| \leq \sum_{i \text{ is good}} \frac{\epsilon}{5} |D_2(i)| \leq \frac{\epsilon}{5},$$

we have

$$\sum_{i \text{ is bad}} (|D_1(i)| + |D_2(i)|) \geq \sum_{i \text{ is bad}} |D_1(i) - D_2(i)| \geq \frac{9}{5}\epsilon.$$

Consequently it must be the case that either $D_1(\text{BAD}) \geq \frac{9}{10}\epsilon$ or $D_2(\text{BAD}) \geq \frac{9}{10}\epsilon$. For the rest of the argument we suppose that $D_1(\text{BAD}) \geq \frac{9}{10}\epsilon$ (by the symmetry of the algorithm, an identical argument to the one we give below but with the roles of D_1 and D_2 flipped throughout handles the other case).

Since $D_1(\text{BAD}) \geq \frac{9}{10}\epsilon$, a simple calculation shows that with probability at least $98/100$ at least one of the $5/\epsilon$ points h_1, \dots, h_m drawn in Line 1 belongs to BAD . For the rest of the argument we suppose that indeed (at least) one of these points is in BAD ; let h_{i^*} be such a point. Now consider the execution of Line 4 when ORACLE_1 is called on h_{i^*} . By Definition 6, whether or not i^* belongs to $S(\epsilon/100, D_1)$, with probability at least $1 - \epsilon/100$ the call to ORACLE_1 either causes $\text{TEST-EQUALITY-UNKNOWN}$ to REJECT in Line 4 (because ORACLE_1 returns UNKNOWN) or it returns a value $v_{1,i^*} \in [1 - \frac{\epsilon}{100}, 1 + \frac{\epsilon}{100}]D_1(i^*)$. We may suppose that it returns a value $v_{1,i^*} \in [1 - \frac{\epsilon}{100}, 1 + \frac{\epsilon}{100}]D_1(i^*)$. Similarly, in the execution of Line 5 when ORACLE_2 is called on h_{i^*} , whether or not i^* belongs to $S(\epsilon/100, D_2)$, with probability at least $1 - \epsilon/100$ the call to ORACLE_2 either causes $\text{TEST-EQUALITY-UNKNOWN}$ to reject in Line 5 or it returns a value $v_{2,i^*} \in [1 - \frac{\epsilon}{100}, 1 + \frac{\epsilon}{100}]D_2(i^*)$. We may suppose that it returns a value $v_{2,i^*} \in [1 - \frac{\epsilon}{100}, 1 + \frac{\epsilon}{100}]D_2(i^*)$. But recalling that $i^* \in \text{BAD}$, an easy calculation shows that the values v_{1,i^*} and v_{2,i^*} must be multiplicatively far enough from each other that the algorithm will output REJECT in Line 7. Thus with overall probability at least $96/100$ the algorithm outputs REJECT . ■

6.2.3 Constructing an approximate EVAL_D simulator using COND_D

In this subsection we show that a COND_D oracle can be used to obtain an approximate EVAL simulator:

Theorem 13 *Let D be any distribution over $[N]$ and let $0 < \epsilon, \delta < 1$. The algorithm $\text{APPROX-EVAL-SIMULATOR}$ has the following properties: It uses*

$$\tilde{O}\left(\frac{(\log N)^5 \cdot (\log(1/\delta))^2}{\epsilon^3}\right)$$

calls to COND_D and it is an (ϵ, δ) -approximate EVAL_D simulator.

A few notes: First, in the proof we give below of Theorem 13 we assume throughout that $0 < \epsilon \leq c$, where c is a small absolute constant. This incurs no loss of generality because if the desired ϵ parameter is in $(c, 1)$ then the parameter can simply be set to $c/2$. We further note that in keeping with our requirement on a COND_D algorithm, the algorithm $\text{APPROX-EVAL-SIMULATOR}$ only ever calls the COND_D oracle on sets S which are either $S = [N]$ or else contain at least one element i that has been returned as the output of an earlier call to COND_D . (To see this, note that Line 6 is the only line when COND_D queries are performed. In the first execution of the outer “For” loop clearly all COND queries are on set $S_0 = [N]$. In subsequent stages the only way a set S_j is formed is if either (i) S_j is set to $\{i^*\}$ in Line 10, in which case clearly i^* was previously received as the response of a $\text{COND}_D(S_{j-1})$ query, or else (ii) a nonzero fraction of elements i_1, \dots, i_m received as responses to $\text{COND}_D(S_{j-1})$ queries belong to S_j (see Line 19).)

A preliminary simplification. Fix a distribution D over $[N]$. Let Z denote $\text{supp}(D)$, i.e. $Z = \{i \in [N] : D(i) > 0\}$. We first claim that in proving Theorem 13 we may assume without loss of generality that no two distinct elements $i, j \in Z$ have $D(i) = D(j)$ – in other words, we shall prove the theorem under this assumption on D , and we claim that this implies the general result. To see this, observe that if Z contains elements $i \neq j$ with $D(i) = D(j)$, then for any arbitrarily small $\xi > 0$ and any arbitrarily large M we can perturb the weights of elements in Z to obtain a distribution D' supported on Z such that (i) no two elements of Z have the same probability under D' , and (ii) for every $S \subseteq [N]$, $S \cap Z \neq \emptyset$ we have $d_{\text{TV}}(D_S, D'_S) \leq \xi/M$. Since the variation distance between D'_S and D_S is at most ξ/M for an arbitrarily small ξ , the variation distance between (the execution of any M -query COND algorithm run on D) and (the execution of any M -query COND algorithm run on D') will be at most ξ . Since ξ can be made arbitrarily small this means that indeed without loss of generality we may work with D' in what follows.

Thus, we henceforth assume that the distribution D has no two elements in $\text{supp}(D)$ with the same weight. For such a distribution we can explicitly describe the set $S^{(\epsilon, D)}$ from Definition 6 that our analysis will deal with. Let $\pi : \{1, \dots, |Z|\} \rightarrow Z$ be the bijection such that $D(\pi(1)) > \dots > D(\pi(|Z|))$ (note that the bijection π is uniquely defined by the assumption that $D(i) \neq D(j)$ for all distinct $i, j \in Z$). Given a value $0 < \tau < 1$ we define the set $L_{\tau, D}$ to be $([N] \setminus Z) \cup \{\pi(s), \dots, \pi(|Z|)\}$ where s is the smallest index in $\{1, \dots, |Z|\}$ such that $\sum_{j=s}^{|Z|} D(\pi(j)) < \tau$ (if $D(\pi(|Z|))$ itself is at least τ then we define $L_{\tau, D} = [N] \setminus Z$). Thus intuitively $L_{\tau, D}$ contains the τ fraction (w.r.t. D) of $[N]$ consisting of the lightest elements. The desired set $S^{(\epsilon, D)}$ is precisely $L_{\epsilon, D}$.

Intuition for the algorithm. The high-level idea of the EVAL_D simulation is the following: Let $i^* \in [N]$ be the input element given to the EVAL_D simulator. The algorithm works in a sequence of stages. Before performing the j -th stage it maintains a set S_{j-1} that contains i^* , and it has a high-accuracy estimate $\hat{D}(S_{j-1})$ of the value of $D(S_{j-1})$. (The initial set S_0 is simply $[N]$ and the initial estimate $\hat{D}(S_0)$ is of course 1.) In the j -th stage the algorithm attempts to construct a subset S_j of S_{j-1} in such a way that (i) $i^* \in S_j$, and (ii) it is possible to obtain a high-accuracy estimate of $D(S_j)/D(S_{j-1})$ (and thus a high-accuracy estimate of $D(S_j)$). If the algorithm cannot construct such a set S_j then it outputs UNKNOWN ; otherwise, after at most (essentially) $O(\log N)$ stages, it reaches a situation where $S_j = \{i^*\}$ and so the high-accuracy estimate of $D(S_j) = D(i^*)$ is the desired value.

A natural first idea towards implementing this high-level plan is simply to split S_{j-1} randomly into two pieces and use one of them as S_j . However this simple approach may not work; for example,

if S_{j-1} has one or more elements which are very heavy compared to i^* , then with a random split it may not be possible to efficiently estimate $D(S_j)/D(S_{j-1})$ as required in (ii) above. Thus we follow a more careful approach which first identifies and removes “heavy” elements from S_{j-1} in each stage.

In more detail, during the j -th stage, the algorithm first performs COND_D queries on the set S_{j-1} to identify a set $H_j \subseteq S_{j-1}$ of “heavy” elements; this set essentially consists of all elements which individually each contribute at least a κ fraction of the total mass $D(S_{j-1})$. (Here κ is a “not-too-small” quantity but it is significantly less than ϵ .) Next, the algorithm performs additional COND_D queries to estimate $D(i^*)/D(S_{j-1})$. If this fraction exceeds $\kappa/20$ then it is straightforward to estimate $D(i^*)/D(S_{j-1})$ to high accuracy, so using $\hat{D}(S_{j-1})$ it is possible to obtain a high-quality estimate of $D(i^*)$ and the algorithm can conclude. However, the typical case is that $D(i^*)/D(S_{j-1}) < \kappa/20$. In this case, the algorithm next estimates $D(H_j)/D(S_{j-1})$. If this is larger than $1 - \epsilon/10$ then the algorithm outputs UNKNOWN (see below for more discussion of this). If $D(H_j)/D(S_{j-1})$ is less than $1 - \epsilon/10$ then $D(S_{j-1} \setminus H_j)/D(S_{j-1}) \geq \epsilon/10$ (and so $D(S_{j-1} \setminus H_j)/D(S_{j-1})$ can be efficiently estimated to high accuracy), but each element k of $S_{j-1} \setminus H_j$ has $D(k)/D(S_{j-1}) \leq \kappa \ll \epsilon/10 \leq D(S_{j-1} \setminus H_j)/D(S_{j-1})$. Thus it must be the case that the weight under D of $S_{j-1} \setminus H_j$ is “spread out” over many “light” elements.

Given that this is the situation, the algorithm next chooses S'_j to be a random subset of $S_{j-1} \setminus (H_j \cup \{i^*\})$, and sets S_j to be $S'_j \cup \{i^*\}$. It can be shown that with high probability (over the random choice of S_j) it will be the case that $D(S_j) \geq \frac{1}{3}D(S_{j-1} \setminus H_j)$ (this relies crucially on the fact that the weight under D of $S_{j-1} \setminus H_j$ is “spread out” over many “light” elements). This makes it possible to efficiently estimate $D(S_j)/D(S_{j-1} \setminus H_j)$; together with the high-accuracy estimate of $D(S_{j-1} \setminus H_j)/D(S_{j-1})$ noted above, and the high-accuracy estimate $\hat{D}(S_{j-1})$ of $D(S_{j-1})$, this means it is possible to efficiently estimate $D(S_j)$ to high accuracy as required for the next stage. (We note that after defining S_j but before proceeding to the next stage, the algorithm actually checks to be sure that S_j contains at least one point that was returned from the $\text{COND}_D(S_{j-1})$ calls made in the past stage. This check ensures that whenever the algorithm calls $\text{COND}_D(S)$ on a set S , it is guaranteed that $D(S) > 0$ as required by our COND_D model. Our analysis shows that doing this check does not affect correctness of the algorithm since with high probability the check always passes.)

Intuition for the analysis. We require some definitions to give the intuition for the analysis establishing correctness. Fix a nonempty subset $S \subseteq [N]$. Let π_S be the bijection mapping $\{1, \dots, |S|\}$ to S in such a way that $D_S(\pi_S(1)) > \dots > D_S(\pi_S(|S|))$, i.e. $\pi_S(1), \dots, \pi_S(|S|)$ is a listing of the elements of S in order from heaviest under D_S to lightest under D_S . Given $j \in S$, we define the S -rank of j , denoted $\text{rank}_S(j)$, to be the value $\sum_{i: D_S(\pi(i)) \leq D_S(j)} D_S(\pi(i))$, i.e. $\text{rank}_S(j)$ is the sum of the weights (under D_S) of all the elements in S that are no heavier than j under D_S . Note that having $i^* \notin L_{\epsilon, N}$ implies that $\text{rank}_{[N]}(i^*) \geq \epsilon$.

We first sketch the argument for correctness. (It is easy to show that the algorithm only outputs FAIL with very small probability so we ignore this possibility below.) Suppose first that $i^* \notin L_{\epsilon, D}$. A key lemma shows that if $i^* \notin L_{\epsilon, D}$ (and hence $\text{rank}_{[N]}(i^*) \geq \epsilon$), then with high probability every set S_{j-1} constructed by the algorithm is such that $\text{rank}_{S_{j-1}}(i^*) \geq \epsilon/2$. (In other words, if i^* is not initially among the ϵ -fraction (under D) of lightest elements, then it never “falls too far” to become part of the $\epsilon/2$ -fraction (under $D_{S_{j-1}}$) of lightest elements for S_{j-1} , for any j .) Given that (whp)

i^* always has $\text{rank}_{S_{j-1}}(i^*) \geq \epsilon/2$, though, then it must be the case that (whp) the procedure does not output UNKNOWN (and hence it must whp output a numerical value). This is because there are only two places where the procedure can output UNKNOWN, in Lines 14 and 19; we consider both cases below.

1. In order for the procedure to output UNKNOWN in Line 14, it must be the case that the elements of H_j – each of which individually has weight at least $\kappa/2$ under $D_{S_{j-1}}$ – collectively have weight at least $1 - 3\epsilon/20$ under $D_{S_{j-1}}$ by Line 13. But i^* has weight at most $3\kappa/40$ under $D_{S_{j-1}}$ (because the procedure did not go to Line 2 in Line 10), and thus i^* would need to be in the bottom $3\epsilon/20$ of the lightest elements, i.e. it would need to have $\text{rank}_{S_{j-1}}(i^*) \leq 3\epsilon/20$; but this contradicts $\text{rank}_{S_{j-1}}(i^*) \geq \epsilon/2$.
2. In order for the procedure to output UNKNOWN in Line 19, it must be the case that all elements i_1, \dots, i_m drawn in Line 6 are not chosen for inclusion in S_j . In order for the algorithm to reach Line 19, though, it must be the case that at least $(\epsilon/10 - \kappa/20)m$ of these draws do not belong to $H_j \cup \{i^*\}$; since these draws do not belong to H_j each one occurs only a small number of times among the m draws, so there must be many distinct values, and hence the probability that none of these distinct values is chosen for inclusion in S'_j is very low.

Thus we have seen that if $i^* \notin L_{\epsilon,D}$, then whp the procedure outputs a numerical value; it remains to show that whp this value is a high-accuracy estimate of $D(i^*)$. However, this follows easily from the fact that we inductively maintain a high-quality estimate of $D(S_{j-1})$ and the fact that the algorithm ultimately constructs its estimate of $\hat{D}(i^*)$ only when it additionally has a high-quality estimate of $D(i^*)/D(S_{j-1})$. This fact also handles the case in which $i^* \in L_{\epsilon,D}$ – in such a case it is allowable for the algorithm to output UNKNOWN, so since the algorithm w.h.p. outputs a high-accuracy estimate when it outputs a numerical value, this means the algorithm performs as required in Case (ii) of Definition 6.

We now sketch the argument for query complexity. We will show that the heavy elements can be identified in each stage using $\text{poly}(\log N, 1/\epsilon)$ queries. Since the algorithm constructs S_j by taking a random subset of S_{j-1} (together with i^*) at each stage, the number of stages is easily bounded by (essentially) $O(\log N)$. Since the final probability estimate for $D(i^*)$ is a product of $O(\log N)$ conditional probabilities, it suffices to estimate each of these conditional probabilities to within a multiplicative factor of $(1 \pm O(\frac{\epsilon}{\log N}))$. We show that each conditional probability estimate can be carried out to this required precision using only $\text{poly}(\log N, 1/\epsilon)$ calls to COND_D ; given this, the overall $\text{poly}(\log N, 1/\epsilon)$ query bound follows straightforwardly.

Now we enter into the actual proof. We begin our analysis with a simple but useful lemma about the “heavy” elements identified in Line 7.

Lemma 16 *With probability at least $1 - \delta/K$, every set H_j that is ever constructed in Line 7 satisfies the following for all $\ell \in S_{j-1}$:*

- (i) *If $D(\ell)/D(S_{j-1}) > \kappa$, then $\ell \in H_j$;*
- (ii) *If $D(\ell)/D(S_{j-1}) < \kappa/2$ then $\ell \notin H_j$.*

Algorithm 10: APPROX-EVAL-SIMULATOR

Input: access to COND_D ; parameters $0 < \epsilon, \delta < 1$; input element $i^* \in [N]$

- 1: Set $S_0 = [N]$ and $\hat{D}(S_0) = 1$. Set $K = 9$. Set $M = \log N + \log(K/\delta) + 1$. Set $\kappa = \Theta(\epsilon/(M^2 \log(M/\delta)))$.
- 2: **for** $j = 1$ to M **do**
- 3: **if** $|S_{j-1}| = 1$ **then**
- 4: return $\hat{D}(S_{j-1})$ (and exit)
- 5: **end if**
- 6: Perform $m = \Theta(\max\{M^2 \log(M/\delta)/(\epsilon^2 \kappa), \log(M/(\delta \kappa))/\kappa^2\})$ COND_D queries on S_{j-1} to obtain points $i_1, \dots, i_m \in S_{j-1}$.
- 7: Let $H_j = \{k \in [N] : k \text{ appears at least } \frac{3}{4} \kappa m \text{ times in the list } i_1, \dots, i_m\}$
- 8: Let $\hat{D}_{S_{j-1}}(i^*)$ denote the fraction of times that i^* appears in i_1, \dots, i_m
- 9: **if** $\hat{D}_{S_{j-1}}(i^*) \geq \frac{\kappa}{20}$ **then**
- 10: Set $S_j = \{i^*\}$, set $\hat{D}(S_j) = \hat{D}_{S_{j-1}}(i^*) \cdot \hat{D}(S_{j-1})$, increment j , and go to Line 2.
- 11: **end if**
- 12: Let $\hat{D}_{S_{j-1}}(H_j)$ denote the fraction of elements among i_1, \dots, i_m that belong to H_j .
- 13: **if** $\hat{D}_{S_{j-1}}(H_j) > 1 - \epsilon/10$ **then**
- 14: return UNKNOWN (and exit)
- 15: **end if**
- 16: Set S'_j to be a uniform random subset of $S_{j-1} \setminus (H_j \cup \{i^*\})$ and set S_j to be $S'_j \cup \{i^*\}$.
- 17: Let $\hat{D}_{S_{j-1}}(S_j)$ denote the fraction of elements among i_1, \dots, i_m that belong to S_j
- 18: **if** $\hat{D}_{S_{j-1}}(S_j) = 0$ **then**
- 19: return UNKNOWN (and exit)
- 20: **end if**
- 21: Set $\hat{D}(S_j) = \hat{D}_{S_{j-1}}(S_j) \cdot \hat{D}(S_{j-1})$
- 22: **end for**
- 23: Output FAIL.

Proof: Fix an iteration j . By Line 7 in the algorithm, a point ℓ is included in H_j if it appears at least $\frac{3}{4}\kappa m$ times among i_1, \dots, i_m (which are the output of COND_D queries on S_{j-1}). For the first item, fix an element ℓ such that $D(\ell)/D(S_{j-1}) > \kappa$. Recall that $m = \Omega(M^2 \log(M/\delta)/(\epsilon^2 \kappa)) = \Omega(\log(MN/\delta)/\kappa)$ (since $M = \Omega(\log(N))$). By a multiplicative Chernoff bound, the probability (over the choice of i_1, \dots, i_m in S_{j-1}) that ℓ appears less than $\frac{3}{4}\kappa m$ times among i_1, \dots, i_m (that is, less than $3/4$ times the lower bound on the expected value) is at most $\delta/(KMN)$ (for an appropriate constant in the setting of m). On the other hand, for each fixed ℓ such that $D(\ell)/D(S_{j-1}) < \kappa/2$, the probability that ℓ appears at least $\frac{3}{4}\kappa m$ times (that is, at least $3/2$ times the upper bound on the expected value) is at most $\delta/(KMN)$ as well. The lemma follows by taking a union bound over all (at most N) points considered above and over all M settings of j , ■

Next we show that with high probability Algorithm APPROX-EVAL-SIMULATOR returns either UNKNOWN or a numerical value (as opposed to outputting FAIL in Line 23):

Lemma 17 *For any D , ϵ, δ and i^* , Algorithm APPROX-EVAL-SIMULATOR outputs FAIL with probability at most δ/K .*

Proof: Fix any element $i \neq i^*$. The probability (taken only over the choice of the random subset in each execution of Line 16) that i is placed in S'_j in each of the first $\log N + \log(K/\delta)$ executions of Line 16 is at most $\frac{\delta}{KN}$. Taking a union bound over all $N - 1$ points $i \neq i^*$, the probability that any point other than i^* remains in S_{j-1} through all of the first $\log N + \log(K/\delta)$ executions of the outer “for” loop is at most $\frac{\delta}{K}$. Assuming that this holds, then in the execution of the outer “for” loop when $j = \log N + \log(K/\delta) + 1$, the algorithm will return $\hat{D}(S_{j-1}) = \hat{D}(i^*)$ in Line 4. ■

For the rest of the analysis it will be helpful for us to define several “desirable” events and show that they all hold with high probability:

1. Let E_1 denote the event that every set H_j that is ever constructed in Line 7 satisfies both properties (i) and (ii) stated in Lemma 16. By Lemma 16 the event E_1 holds with probability at least $1 - \delta/K$.
2. Let E_2 denote the event that in every execution of Line 9, the estimate $\hat{D}_{S_{j-1}}(i^*)$ is within an additive $\pm \frac{\kappa}{40}$ of the true value of $D(i^*)/D(S_{j-1})$. By the choice of m in Line 6 (i.e., using $m = \Omega(\log(M/\delta)/\kappa^2)$), an additive Chernoff bound, and a union bound over all iterations, the event E_2 holds with probability at least $1 - \delta/K$.
3. Let E_3 denote the event that if Line 10 is executed, the resulting value $\hat{D}_{S_{j-1}}(i^*)$ lies in $[1 - \frac{\epsilon}{2M}, 1 + \frac{\epsilon}{2M}]D(i^*)/D(S_{j-1})$. Assuming that event E_2 holds, if Line 10 is reached then the true value of $D(i^*)/D(S_{j-1})$ must be at least $\kappa/40$, and consequently a multiplicative Chernoff bound and the choice of m (i.e. using $m = \Omega(M^2 \log(M/\delta)/(\epsilon^2 \kappa))$) together imply that $\hat{D}_{S_{j-1}}(i^*)$ lies in $[1 - \frac{\epsilon}{2M}, 1 + \frac{\epsilon}{2M}]D(i^*)/D(S_{j-1})$ except with failure probability at most δ/K .
4. Let E_4 denote the event that in every execution of Line 12, the estimate $\hat{D}_{S_{j-1}}(H_j)$ is within an additive error of $\pm \frac{\epsilon}{20}$ from the true value of $D(H_j)/D(S_{j-1})$. By the choice of m in Line 6 (i.e., using $m = \Omega(\log(M/\delta)/\epsilon^2)$) and an additive Chernoff bound, the event E_4 holds with probability at least $1 - \delta/K$.

The above arguments show that E_1, E_2, E_3 and E_4 all hold with probability at least $1 - 4\delta/K$.

Let E_5 denote the event that in every execution of Line 16, the set S'_j which is drawn satisfies $D(S'_j)/D(S_{j-1} \setminus (H_j \cup \{i^*\})) \geq 1/3$. The following lemma says that conditioned on E_1 through E_4 all holding, event E_5 holds with high probability:

Lemma 18 *Conditioned on E_1 through E_4 the probability that E_5 holds is at least $1 - \delta/K$.*

Proof: Fix a value of j and consider the j -th iteration of Line 16. Since events E_2 and E_4 hold, it must be the case that $D(S_{j-1} \setminus (H_j \cup \{i^*\}))/D(S_{j-1}) \geq \epsilon/40$. Since event E_1 holds, it must be the case that every $i \in (S_{j-1} \setminus (H_j \cup \{i^*\}))$ has $D(i)/D(S_{j-1}) \leq \kappa$. Now since S'_j is chosen by independently including each element of $S_{j-1} \setminus (H_j \cup \{i^*\})$ with probability $1/2$, we can apply the first part of Corollary 3 and get

$$\Pr \left[D(S'_j) < \frac{1}{3} D(S_{j-1} \setminus (H_j \cup \{i^*\})) \right] \leq \epsilon^{-4\epsilon/(40 \cdot 9 \cdot 4\kappa)} < \frac{\delta}{KM},$$

where the last inequality follows by the setting of $\kappa = \Theta(\epsilon/(M^2 \log(1/\delta)))$. ■

Thus we have established that E_1 through E_5 all hold with probability at least $1 - 5\delta/K$.

Next, let E_6 denote the event that the algorithm never returns UNKNOWN and exits in Line 19. Our next lemma shows that conditioned on events E_1 through E_5 , the probability of E_6 is at least $1 - \delta/K$:

Lemma 19 *Conditioned on E_1 through E_5 the probability that E_6 holds is at least $1 - \delta/K$.*

Proof: Fix any iteration j of the outer “For” loop. In order for the algorithm to reach Line 18 in this iteration, it must be the case (by Lines 9 and 13) that at least $(\epsilon/10 - \kappa/20)m > (\epsilon/20)m$ points in i_1, \dots, i_m do not belong to $H_j \cup \{i^*\}$. Since each point not in H_j appears at most $\frac{3}{4}\kappa m$ times in the list i_1, \dots, i_m , there must be at least $\frac{\epsilon}{15\kappa}$ distinct such values. Hence the probability that none of these values is selected to belong to S'_j is at most $1/2^{\epsilon/(15\kappa)} < \delta/(KM)$. A union bound over all (at most M) values of j gives that the probability the algorithm ever returns UNKNOWN and exits in Line 19 is at most δ/M , so the lemma is proved. ■

Now let E_7 denote the event that in every execution of Line 17, the estimate $\hat{D}_{S_{j-1}}(S_j)$ lies in $[1 - \frac{\epsilon}{2M}, 1 + \frac{\epsilon}{2M}]D(S_j)/D(S_{j-1})$. The following lemma says that conditioned on E_1 through E_5 , event E_7 holds with probability at least $1 - \delta/K$:

Lemma 20 *Conditioned on E_1 through E_5 , the probability that E_7 holds is at least $1 - \delta/K$.*

Proof: Fix a value of j and consider the j -th iteration of Line 16. The expected value of $\hat{D}_{S_{j-1}}(S_j)$ is precisely

$$\frac{D(S_j)}{D(S_{j-1})} = \frac{D(S_j)}{D(S_{j-1} \setminus (H_j \cup \{i^*\}))} \cdot \frac{D(S_{j-1} \setminus (H_j \cup \{i^*\}))}{D(S_{j-1})}. \quad (59)$$

Since events E_2 and E_4 hold we have that $\frac{D(S_{j-1} \setminus (H_j \cup \{i^*\}))}{D(S_{j-1})} \geq \epsilon/40$, and since event E_5 holds we have that $\frac{D(S_j)}{D(S_{j-1} \setminus (H_j \cup \{i^*\}))} \geq 1/3$ (note that $D(S_j) \geq D(S'_j)$). Thus we have that (59) is at least $\epsilon/120$. Recalling the value of m (i.e., using $m = \Omega(M^2 \log(M/\delta)/\epsilon^2 \kappa) = \Omega(M^2 \log(KM/\delta)/\epsilon^3)$) a multiplicative Chernoff bound gives that indeed $\hat{D}_{S_{j-1}}(S_j) \in [1 - \frac{\epsilon}{2M}, 1 + \frac{\epsilon}{2M}]D(S_j)/D(S_{j-1})$ with failure probability at most $\delta/(KM)$. A union bound over all M possible values of j finishes the proof. ■

At this point we have established that events E_1 through E_7 all hold with probability at least $1 - 7\delta/K$.

We can now argue that each estimate $\hat{D}(S_j)$ is indeed a high-accuracy estimate of the true value $D(S_j)$:

Lemma 21 *With probability at least $1 - 7\delta/K$ each estimate $\hat{D}(S_j)$ constructed by APPROX-EVAL-SIMULATOR lies in $[(1 - \frac{\epsilon}{2M})^j, (1 + \frac{\epsilon}{2M})^j]D(S_j)$.*

Proof: We prove the lemma by showing that if all events E_1 through E_7 hold, then the following claim (denoted $(*)$) holds: each estimate $\hat{D}(S_j)$ constructed by APPROX-EVAL-SIMULATOR lies in $[(1 - \frac{\epsilon}{2M})^j, (1 + \frac{\epsilon}{2M})^j]D(S_j)$. Thus for the rest of the proof we assume that indeed all events E_1 through E_7 hold.

The claim $(*)$ is clearly true for $j = 0$. We prove $(*)$ by induction on j assuming it holds for $j - 1$. The only places in the algorithm where $\hat{D}(S_j)$ may be set are Lines 10 and 21. If $\hat{D}(S_j)$ is set in Line 21 then $(*)$ follows from the inductive claim for $j - 1$ and Lemma 20. If $\hat{D}(S_j)$ is set in Line 10, then $(*)$ follows from the inductive claim for $j - 1$ and the fact that event E_3 holds. This concludes the proof of the lemma. ■

Finally, we require the following crucial lemma which establishes that if $i^* \notin L_{\epsilon, N}$ (and hence the initial rank $\text{rank}_{[N]}(i^*)$ is at least ϵ), then with very high probability the rank of i^* never becomes too low during the execution of the algorithm:

Lemma 22 *Suppose $i^* \notin L_{\epsilon, N}$. Then with probability at least $1 - \delta/K$, every set S_{j-1} constructed by the algorithm has $\text{rank}_{S_{j-1}}(i^*) \geq \epsilon/2$.*

We prove Lemma 22 in Section 6.2.4 below.

With these pieces in place we are ready to prove Theorem 13.

Proof of Theorem 13: It is straightforward to verify that algorithm APPROX-EVAL-SIMULATOR has the claimed query complexity. We now argue that APPROX-EVAL-SIMULATOR meets the two requirements (i) and (ii) of Definition 6. Throughout the discussion below we assume that all the “favorable events” in the above analysis (i.e. events E_1 through E_7 , Lemma 17, and Lemma 22) indeed hold as desired (incurring an overall failure probability of at most δ).

Suppose first that $i^* \notin L_{\epsilon, D}$. By Lemma 22 it must be the case that the algorithm does not return UNKNOWN in Line 14. (This is because in order to reach Line 14 it would need to be the case that $D(i^*)/D(S_{j-1}) \leq 3\kappa/40$ (so the algorithm does not instead go to Line 22 in Line 10), but

since by Lemma 16 every element k in H_j has $D(k)/D(S_{j-1}) \geq \kappa/2$, this means that i^* does not belong to H_j . In order to reach Line 14, by event E_4 we must have $D(H_j)/D(S_{j-1}) \geq 1 - 3\epsilon/20$. Since every element of H_j has more mass under D (at least $\kappa/2$) than i^* (which has at most $3\kappa/40$), this implies that $\text{rank}_{S_{j-1}}(i^*) \leq 3\epsilon/20$. This contradicts Lemma 22.) And by Lemma 19 it must be the case that the algorithm does not return UNKNOWN in Line 19. Thus the algorithm terminates by returning an estimate $\hat{D}(S_j) = \hat{D}(i^*)$ which, by Lemma 21, lies in $[(1 - \frac{\epsilon}{2M})^j, (1 + \frac{\epsilon}{2M})^j]D(i^*)$. Since $j \leq M$ this estimate lies in $[1 - \epsilon, 1 + \epsilon]D(i^*)$ as required.

Now suppose that $i^* \in L_{\epsilon,D}$. By Lemma 17 we may assume that the algorithm either outputs UNKNOWN or a numerical value. As above, Lemma 21 implies that if the algorithm outputs a numerical value then the value lies in $[1 - \epsilon, 1 + \epsilon]D(i^*)$ as desired. This concludes the proof of Theorem 13. ■

6.2.4 Proof of Lemma 22.

The key to proving Lemma 22 will be proving the next lemma. (In the following, for S a set of real numbers we write $\text{sum}(S)$ to denote $\sum_{\alpha \in S} \alpha$.)

Lemma 23 *Fix $0 < \epsilon \leq c$. Set $\kappa = \Theta(\epsilon/(M^2 \log(1/\delta)))$. Let $T = \{\alpha_1, \dots, \alpha_n\}$ be a set of values $\alpha_1 < \dots < \alpha_n$ such that $\text{sum}(T) = 1$. Fix $\ell \in [n]$ and let $T_L = \{\alpha_1, \dots, \alpha_\ell\}$ and let $T_R = \{\alpha_{\ell+1}, \dots, \alpha_n\}$, so $T_L \cup T_R = T$. Assume that $\text{sum}(T_L) \geq \epsilon/2$ and that $\alpha_\ell \leq \kappa/10$.*

Fix H to be any subset of T satisfying the following two properties: (i) H includes every α_j such that $\alpha_j \geq \kappa$; and (ii) H includes no α_j such that $\alpha_j < \kappa/2$. (Note that consequently H does not intersect T_L .)

Let T' be a subset of $(T \setminus (H \cup \{\alpha_\ell\}))$ selected uniformly at random. Let $T'_L = T' \cap T_L$ and let $T'_R = T' \cap T_R$.

Then we have the following:

1. *If $\text{sum}(T_L) > 20\epsilon$, then with probability at least $1 - \delta/M$ (over the random choice of T') it holds that*

$$\frac{\text{sum}(T'_L \cup \{\alpha_\ell\})}{\text{sum}(T' \cup \{\alpha_\ell\})} \geq 9\epsilon;$$

2. *If $\epsilon/2 \leq \text{sum}(T_L) < 20\epsilon$, then with probability at least $1 - \delta/M$ (over the random choice of T') it holds that*

$$\frac{\text{sum}(T'_L \cup \{\alpha_\ell\})}{\text{sum}(T' \cup \{\alpha_\ell\})} \geq \text{sum}(T_L) (1 - \rho),$$

where $\rho = \frac{\ln 2}{M}$.

Proof of Lemma 22 using Lemma 23: We apply Lemma 23 repeatedly at each iteration j of the outer “For” loop. The set H of Lemma 23 corresponds to the set H_j of “heavy” elements that are removed at a given iteration, the set of values T corresponds to the values $D(i)/D(S_{j-1})$

for $i \in S_{j-1}$, and the element α_ℓ of Lemma 23 corresponds to $D(i^*)/D(S_{j-1})$. The value $\text{sum}(T_L)$ corresponds to $\text{rank}_{S_{j-1}}(i^*)$ and the value

$$\frac{\text{sum}(T'_L \cup \{\alpha_\ell\})}{\text{sum}(T' \cup \{\alpha_\ell\})}$$

corresponds to $\text{rank}_{S_j}(i^*)$. Observe that since $i^* \notin L_{\epsilon, N}$ we know that initially $\text{rank}_{[N]}(i^*) \geq \epsilon$, which means that the first time we apply Lemma 23 (with $T = \{D(i) : i \in [N]\}$) we have $\text{sum}(T_L) \geq \epsilon$.

By Lemma 23 the probability of failure in any of the (at most M) iterations is at most δ/K , so we assume that there is never a failure. Consequently for all j we have that if $\text{rank}_{S_{j-1}}(i^*) \geq 20\epsilon$ then $\text{rank}_{S_j}(i^*) \geq 9\epsilon$, and if $\epsilon/2 \leq \text{rank}_{S_{j-1}}(i^*) < 20\epsilon$ then $\text{rank}_{S_j}(i^*) \geq \text{rank}_{S_{j-1}}(i^*) \cdot (1 - \rho)$. Since $\text{rank}_{S_0}(i^*) \geq \epsilon$, it follows that for all $j \leq M$ we have $\text{rank}_{S_j}(i^*) \geq \epsilon \cdot (1 - \rho)^M > \epsilon/2$. ■

Proof of Lemma 23. We begin with the following claim:

Claim 24 *With probability at least $1 - \delta/(2M)$ (over the random choice of T') it holds that $\text{sum}(T'_L) \geq \frac{1}{2} \cdot \text{sum}(T_L) \cdot (1 - \rho/2)$.*

Proof: Recall from the setup that every element $\alpha_i \in T_L$ satisfies $\alpha_i \leq \kappa/10$, and $\text{sum}(T_L) \geq \epsilon/2$. Also recall that $\kappa = \Theta(\epsilon/(M^2 \log(1/\delta)))$ and that $\rho = \frac{\ln 2}{M}$, so that $\rho^2 \epsilon / (6\kappa) \geq \ln(2M/\delta)$. The claim follows by applying the first part of Corollary 3 (with $\gamma = \rho/2$). ■

Part (1) of Lemma 23 is an immediate consequence of Claim 24, since in part (1) we have

$$\frac{\text{sum}(T'_L \cup \{\alpha_\ell\})}{\text{sum}(T' \cup \{\alpha_\ell\})} \geq \text{sum}(T'_L) \geq \frac{1}{2} \cdot \text{sum}(T_L) \cdot \left(1 - \frac{\rho}{2}\right) \geq \frac{1}{2} \cdot 20\epsilon \cdot \left(1 - \frac{\rho}{2}\right) \geq 9\epsilon.$$

It remains to prove Part (2) of the lemma. We will do this using the following claim:

Claim 25 *Suppose $\epsilon/2 \leq \text{sum}(T_L) \leq 20\epsilon$. Then with probability at least $1 - \delta/(2M)$ (over the random choice of T') it holds that $\text{sum}(T'_R) \leq \frac{1}{2} \text{sum}(T_R) \cdot (1 + \rho/2)$.*

Proof: Observe first that that $\alpha_i < \kappa$ for each $\alpha_i \in T_R \setminus H$. We consider two cases.

If $\text{sum}(T_R \setminus H) \geq 4\epsilon$, then we apply the first part of Corollary 3 to the α_i 's in $T_R \setminus H$ and get that

$$\Pr \left[\text{sum}(T'_R) > \frac{1}{2} \text{sum}(T_R) \cdot (1 + \rho/2) \right] \leq \Pr \left[\text{sum}(T'_R) > \frac{1}{2} \text{sum}(T_R \setminus H) \cdot (1 + \rho/2) \right] \tag{60}$$

$$< \exp(-\rho^2 \text{sum}(T_R \setminus H) / 24\kappa) \tag{60}$$

$$\leq \exp(-\rho^2 \epsilon / (6\kappa)) \leq \frac{\delta}{(2M)} \tag{61}$$

(recall from the proof of Claim 24 that $\rho^2 \epsilon / (6\kappa) \geq \ln(2M/\delta)$).

If $\text{sum}(T_R \setminus H) < 4\epsilon$, (so that the expected value of $\text{sum}(T'_R)$ is less than 2ϵ) then we can apply the second part of Corollary 3 as we explain next. Observe that by the premise of the lemma,

$\text{sum}(T_R) \geq 1 - 20\epsilon$ which is at least $1/2$ (recalling that ϵ is at most a small absolute constant c). Consequently, the event “ $\text{sum}(T'_R) \geq \frac{1}{2} \cdot \text{sum}(T_R) \cdot (1 + \rho/2)$ ” implies the event “ $\text{sum}(T'_R) \geq \frac{1}{4}$ ”, and by applying the second part of Corollary 3 we get

$$\Pr \left[\text{sum}(T'_R) > \frac{1}{2} \text{sum}(T_R) \cdot (1 + \rho/2) \right] \leq \Pr \left[\text{sum}(T'_R) > \frac{1}{4} \right] < 2^{-1/4\kappa} < \frac{\delta}{2}, \quad (62)$$

as required. ■

Now we can prove Lemma 23. Using Claims 24 and 25 we have that with probability at least $1 - \delta/M$,

$$\text{sum}(T'_L) \geq \frac{1}{2} \cdot \text{sum}(T_L) \cdot (1 - \rho/2) \quad \text{and} \quad \text{sum}(T'_R) \leq \frac{1}{2} \text{sum}(T_R) \cdot (1 + \rho/2);$$

we assume that both these inequalities hold going forth. Since

$$\frac{\text{sum}(T'_L \cup \{\alpha_\ell\})}{\text{sum}(T' \cup \{\alpha_\ell\})} = \frac{\text{sum}(T'_L) + \alpha_\ell}{\text{sum}(T') + \alpha_\ell} > \frac{\text{sum}(T'_L)}{\text{sum}(T')},$$

it is sufficient to show that $\frac{\text{sum}(T'_L)}{\text{sum}(T')} \geq \text{sum}(T_L)(1 - \rho)$; we now show this. As $\text{sum}(T') = \text{sum}(T'_L) + \text{sum}(T'_R)$,

$$\begin{aligned} \frac{\text{sum}(T'_L)}{\text{sum}(T')} &= \frac{\text{sum}(T'_L)}{\text{sum}(T'_L) + \text{sum}(T'_R)} = \frac{1}{1 + \frac{\text{sum}(T'_R)}{\text{sum}(T'_L)}} \\ &\geq \frac{1}{1 + \frac{(1/2) \cdot \text{sum}(T_R) \cdot (1 + \rho/2)}{(1/2) \cdot \text{sum}(T_L) \cdot (1 - \rho/2)}} \\ &= \frac{\text{sum}(T_L) \cdot (1 - \rho/2)}{\text{sum}(T_L) \cdot (1 - \rho/2) + \text{sum}(T_R) \cdot (1 + \rho/2)} \\ &\geq \frac{\text{sum}(T_L) \cdot (1 - \rho/2)}{\text{sum}(T_L) \cdot (1 + \rho/2) + \text{sum}(T_R) \cdot (1 + \rho/2)} \\ &= \text{sum}(T_L) \cdot \frac{1 - \rho/2}{1 + \rho/2} > \text{sum}(T_L) \cdot (1 - \rho). \end{aligned}$$

This concludes the proof of Lemma 23. ■

7 An algorithm for estimating the distance to uniformity

In this section we describe an algorithm that estimates the distance between a distribution D and the uniform distribution \mathcal{U} by performing $\text{poly}(1/\epsilon)$ PCOND (and SAMP) queries. We start by giving a high level description of the algorithm.

By the definition of the variation distance (and the uniform distribution),

$$d_{\text{TV}}(D, \mathcal{U}) = \sum_{i: D(i) < 1/N} \left(\frac{1}{N} - D(i) \right). \quad (63)$$

We define the following function over $[N]$:

$$\psi^D(i) = (1 - N \cdot D(i)) \text{ for } D(i) < \frac{1}{N}, \text{ and } \psi^D(i) = 0 \text{ for } D(i) \geq \frac{1}{N}. \quad (64)$$

Observe that $\psi^D(i) \in [0, 1]$ for every $i \in [N]$ and

$$d_{\text{TV}}(D, \mathcal{U}) = \frac{1}{N} \sum_{i=1}^N \psi^D(i). \quad (65)$$

Thus $d_{\text{TV}}(D, \mathcal{U})$ can be viewed as an average value of a function whose range is in $[0, 1]$. Since D is fixed throughout this subsection, we shall use the shorthand $\psi(i)$ instead of $\psi^D(i)$. Suppose we were able to compute $\psi(i)$ exactly for any i of our choice. Then we could obtain an estimate \hat{d} of $d_{\text{TV}}(D, \mathcal{U})$ to within an additive error of $\epsilon/2$ by simply selecting $\Theta(1/\epsilon^2)$ points in $[N]$ uniformly at random and setting \hat{d} to be the average value of $\psi(\cdot)$ on the sampled points. By an additive Chernoff bound (for an appropriate constant in the $\Theta(\cdot)$ notation), with high constant probability the estimate \hat{d} would deviate by at most $\epsilon/2$ from $d_{\text{TV}}(D, \mathcal{U})$.

Suppose next that instead of being able to compute $\psi(i)$ exactly, we were able to compute an estimate $\hat{\psi}(i)$ such that $|\hat{\psi}(i) - \psi(i)| \leq \epsilon/2$. By using $\hat{\psi}(i)$ instead of $\psi(i)$ for each of the $\Theta(1/\epsilon^2)$ sampled points we would incur an additional additive error of at most $\epsilon/2$. Observe first that for i such that $D(i) \leq \epsilon/(2N)$ we have that $\psi(i) \geq 1 - \epsilon/2$, so the estimate $\hat{\psi}(i) = 1$ meets our requirements. Similarly, for i such that $D(i) \geq 1/N$, any estimate $\hat{\psi}(i) \in [0, \epsilon/2]$ can be used. Finally, for i such that $D(i) \in [\epsilon/(2N), 1/N]$, if we can obtain an estimate $\hat{D}(i)$ such that $\hat{D}(i) \in [1 - \epsilon/2, 1 + \epsilon/2]D(i)$, then we can use $\hat{\psi}(i) = N \cdot \hat{D}(i)$.

In order to obtain such estimates $\hat{\psi}(i)$, we shall be interested in finding a *reference point* x . Namely, we shall be interested in finding a pair $(x, \hat{D}(x))$ such that $\hat{D}(x) \in [1 - \epsilon/c, 1 + \epsilon/c]D(x)$ for some sufficiently large constant c , and such that $D(x) = \Omega(\epsilon/N)$ and $D(x) = O(1/(\epsilon N))$. In Subsection 7.1 we describe a procedure for finding such a reference point. More precisely, the procedure is required to find such a reference point (with high constant probability) only under a certain condition on D . It is not hard to verify (and we show this subsequently), that if this condition is not met, then $d_{\text{TV}}(D, \mathcal{U})$ is very close to 1. In order to state the lemma we introduce the following notation. For $\gamma \in [0, 1]$, let

$$H_\gamma^D \stackrel{\text{def}}{=} \left\{ i : D(i) \geq \frac{1}{\gamma N} \right\}. \quad (66)$$

Lemma 26 *Given an input parameter $\kappa \in (0, 1/4]$ as well as SAMP and PCOND query access to a distribution D , the procedure FIND-REFERENCE (Algorithm 12) either returns a pair $(x, \hat{D}(x))$ where $x \in [N]$ and $\hat{D}(x) \in [0, 1]$ or returns No-Pair. The procedure satisfies the following:*

1. *If $D(H_\kappa^D) \leq 1 - \kappa$, then with probability at least $9/10$, the procedure returns a pair $(x, \hat{D}(x))$ such that $\hat{D}(x) \in [1 - 2\kappa, 1 + 3\kappa]D(x)$ and $D(x) \in [\frac{\kappa}{8}, \frac{4}{\kappa}] \cdot \frac{1}{N}$.*
2. *If $D(H_\kappa^D) > 1 - \kappa$, then with probability at least $9/10$, the procedure either returns No-Pair or it returns a pair $(x, \hat{D}(x))$ such that $\hat{D}(x) \in [1 - 2\kappa, 1 + 3\kappa]D(x)$ and $D(x) \in [\frac{\kappa}{8}, \frac{4}{\kappa}] \cdot \frac{1}{N}$.*

The procedure performs $\tilde{O}(1/\kappa^{20})$ PCOND and SAMP queries.

Once we have a reference point x we can use it to obtain an estimate $\hat{\psi}(i)$ for any i of our choice, using the procedure COMPARE, whose properties are stated in Lemma 2 (see Subsection 3.1).

Algorithm 11: Estimating the Distance to Uniformity

Input: PCOND and SAMP query access to a distribution D and a parameter $\epsilon \in [0, 1]$.

1. Call the procedure FIND-REFERENCE (Algorithm 12) with κ set to $\epsilon/8$. If it returns **No-Pair**, then output $\hat{d} = 1$ as the estimate for the distance to uniformity. Otherwise, let $(x, \hat{D}(x))$ be its output.
 2. Select a sample S of $\Theta(1/\epsilon^2)$ points uniformly.
 3. Let $K = \max \left\{ \frac{2/N}{\hat{D}(x)}, \frac{\hat{D}(x)}{\epsilon/(4N)} \right\}$.
 4. For each point $y \in S$:
 - (a) Call COMPARE $\left(\{x\}, \{y\}, \epsilon/2, K, \frac{1}{10|S|} \right)$.
 - (b) If COMPARE returns **High** or it returns a value $\rho(y)$ such that $\rho(y) \cdot \hat{D}(x) \geq 1/N$, then set $\hat{\psi}(y) = 0$;
 - (c) Else, if COMPARE returns **Low** or it returns a value $\rho(y)$ such that $\rho(y) \cdot \hat{D}(x) \leq \epsilon/4N$, then set $\hat{\psi}(y) = 1$;
 - (d) Else set $\hat{\psi}(y) = N \cdot \rho(y) \cdot \hat{D}(x)$.
 5. Output $\hat{d} = \frac{1}{|S|} \sum_{y \in S} \hat{\psi}(y)$.
-

Theorem 14 *With probability at least $2/3$, the estimate \hat{d} returned by Algorithm 11 satisfies: $\hat{d} = d_{\text{TV}}(D, \mathcal{U}) \pm O(\kappa)$. The number of queries performed by the algorithm is $\tilde{O}(1/\epsilon^{20})$.*

Proof: In what follows we shall use the shorthand H_γ instead of H_γ^D . Let E_0 denote the event that the procedure FIND-REFERENCE (Algorithm 12) obeys the requirements in Lemma 26, where by Lemma 26 the event E_0 holds with probability at least $9/10$. Conditioned on E_0 , the algorithm outputs $\hat{d} = 1$ right after calling the procedure (because the procedure returns **No-Pair**) only when $D(H_\kappa) > 1 - \kappa = 1 - \epsilon/8$. We claim that in this case $d_{\text{TV}}(D, \mathcal{U}) \geq 1 - 2\epsilon/8 = 1 - \epsilon/4$. To verify this, observe that

$$d_{\text{TV}}(D, \mathcal{U}) = \sum_{i: D(i) > 1/N} \left(D(i) - \frac{1}{N} \right) \geq \sum_{i \in H_\kappa} \left(D(i) - \frac{1}{N} \right) = D(H_\kappa) - \frac{|H_\kappa|}{N} \geq D(H_\kappa) - \kappa. \quad (67)$$

Thus, in this case the estimate \hat{d} is as required.

We turn to the case in which $D(H_\kappa) \leq 1 - \kappa$ and the procedure FIND-REFERENCE returns a pair $(x, \hat{D}(x))$ such that $\hat{D}(x) \in [1 - 2\kappa, 1 + 3\kappa]D(x)$ and $D(x) \in \left[\frac{\kappa}{8}, \frac{4}{\kappa} \right] \cdot \frac{1}{N}$.

We start by defining two more “desirable” events, which hold (simultaneously) with high constant probability, and then show that conditioned on these events holding (as well as E_0), the output of the algorithm is as required. Let E_1 be the event that the sample S satisfies

$$\left| \frac{1}{|S|} \sum_{y \in S} \psi(y) - d_{\text{TV}}(D, \mathcal{U}) \right| \leq \epsilon/2. \quad (68)$$

By an additive Chernoff bound, the event E_1 holds with probability at least $9/10$.

Next, let E_2 be the event that all calls to the procedure COMPARE return answers as specified in Lemma 2. Since COMPARE is called $|S|$ times, and for each call the probability that it does not return an answer as specified in the lemma is at most $1/(10|S|)$, by the union bound the probability that E_2 holds is at least $9/10$.

From this point on assume events E_0 , E_1 and E_2 all occur, which holds with probability at least $1 - 3/10 \geq 2/3$. Since E_2 holds, we get the following.

1. When COMPARE returns High for $y \in S$ (so that $\hat{\psi}(y)$ is set to 0) we have that

$$D(y) > K \cdot D(x) \geq \frac{2/N}{\widehat{D}(x)} \cdot D(x) > \frac{1}{N}, \quad (69)$$

implying that $\hat{\psi}(y) = \psi(y)$.

2. When COMPARE returns Low for $y \in S$ (so that $\hat{\psi}(y)$ is set to 1) we have that

$$D(y) < \frac{D(x)}{K} \leq \frac{D(x)}{\widehat{D}(x)/(\epsilon/4N)} \leq \frac{\epsilon}{2N}, \quad (70)$$

implying that $\hat{\psi}(y) \leq \psi(y) + \epsilon/2$ (and clearly $\psi(y) \leq \hat{\psi}(y)$).

3. When COMPARE returns a value $\rho(y)$ it holds that $\rho(y) \in [1 - \kappa, 1 + \kappa](D(y)/D(x))$, so that $\rho(y) \cdot \widehat{D}(x) \in [(1 - \kappa)^2, (1 + \kappa)^2]D(y)$. Since $\kappa = \epsilon/8$, if $\rho(y) \cdot \widehat{D}(x) \geq 1/N$ (so that $\hat{\psi}(y)$ is set to 0), then $\psi(y) < \epsilon/2$, if $\rho(y) \cdot \widehat{D}(x) \leq \epsilon/4N$ (so that $\hat{\psi}(y)$ is set to 1), then $\psi(y) \geq 1 - \epsilon/2$, and otherwise $|\hat{\psi}(y) - \psi(y)| \leq \epsilon/2$.

It follows that

$$\hat{d} = \frac{1}{|S|} \sum_{y \in S} \hat{\psi}(y) \in \left[\frac{1}{|S|} \sum_{y \in S} \psi(y) - \epsilon/2, \frac{1}{|S|} \sum_{y \in S} \psi(y) + \epsilon/2 \right] \subseteq [d_{\text{TV}}(D, \mathcal{U}) - \epsilon, d_{\text{TV}}(D, \mathcal{U}) + \epsilon] \quad (71)$$

as required.

The number of queries performed by the algorithm is the number of queries performed by the procedure FIND-REFERENCE, which is $\tilde{O}(1/\epsilon^{20})$, plus $\Theta(1/\epsilon^2)$ times the number of queries performed in each call to COMPARE. The procedure COMPARE is called with the parameter K , which is bounded by $O(1/\epsilon^2)$, the parameter η , which is $\Omega(\epsilon)$, and δ , which is $\Omega(1/\epsilon^2)$. By Lemma 2, the number of queries performed in each call to COMPARE is $O(\log(1/\epsilon)/\epsilon^4)$. The total number of queries performed is hence $\tilde{O}(1/\epsilon^{20})$. ■

7.1 Finding a reference point

In this subsection we prove Lemma 26. We start by giving the high-level idea behind the procedure. For a point $x \in [N]$ and $\gamma \in [0, 1]$, let $U_\gamma^D(x)$ be as defined in Equation (11). Since D is fixed throughout this subsection, we shall use the shorthand $U_\gamma(x)$ instead of $U_\gamma^D(x)$. Recall that κ is a parameter given to the procedure. Assume we had a point x for which $D(U_\kappa(x)) \geq \kappa^{d_1}$ and $|U_\kappa(x)| \geq \kappa^{d_2}N$ for some constants d_1 and d_2 (so that necessarily $D(x) = \Omega(\kappa^{d_1}/N)$ and $D(x) = O(1/(\kappa^{d_2}N))$). It is not hard to verify (and we show this in detail subsequently), that if $D(H) \leq 1 - \kappa$, then a sample of size $\Theta(1/\text{poly}(\kappa))$ distributed according to D will contain such a point x with high constant probability. Now suppose that we could obtain an estimate \hat{w} of $D(U_\kappa(x))$ such that $\hat{w} \in [1 - \kappa, 1 + \kappa]D(U_\kappa(x))$ and an estimate \hat{u} of $|U_\kappa(x)|$ such that $\hat{u} \in [1 - \kappa, 1 + \kappa]|U_\kappa(x)|$. By the definition of $U_\kappa(x)$ we have that $(\hat{w}/\hat{u}) \in [1 - O(\kappa), 1 + O(\kappa)]D(x)$.

Obtaining good estimates of $D(U_\kappa(x))$ and $|U_\kappa(x)|$ (for x such that both $|U_\kappa(x)|$ and $D(U_\kappa(x))$ are sufficiently large) might be infeasible. This is due to the possible existence of many points y for which $D(y)$ is very close to $(1 + \kappa)D(x)$ or $D(x)/(1 + \kappa)$ which define the boundaries of the set $U_\kappa(x)$. For such points it is not possible to efficiently distinguish between those among them that belong to $U_\kappa(x)$ (so that they are within the borders of the set) and those that do not belong to $U_\kappa(x)$ (so that they are just outside the borders of the set). However, for our purposes it suffices to estimate the weight and size of *some* set $U_\alpha(x)$ such that $\alpha \geq \kappa$ (so that $U_\kappa(x) \subseteq U_\alpha(x)$) and α is not much larger than κ (e.g., $\alpha \leq 2\kappa$). To this end we can apply Procedure ESTIMATE-NEIGHBORHOOD (see Subsection 3.2), which (conditioned on $D(U_\kappa(x))$ being above a certain threshold), returns a pair $(\hat{w}(x), \alpha)$ such that $\hat{w}(x)$ is a good estimate of $D(U_\alpha(x))$. Furthermore, α is such that for α' slightly larger than α , the weight of $U_{\alpha'}(x) \setminus U_\alpha(x)$ is small, allowing us to obtain also a good estimate $\hat{\mu}(x)$ of $|U_\alpha(x)|/N$.

Proof of Lemma 26: We first introduce the following notation.

$$L \stackrel{\text{def}}{=} \left\{ i : D(i) < \frac{\kappa}{2N} \right\}, \quad M \stackrel{\text{def}}{=} \left\{ i : \frac{\kappa}{2N} \leq D(i) < \frac{1}{\kappa N} \right\}. \quad (72)$$

Let $H = H_\kappa^D$ where H_κ^D is as defined in Equation (66). Observe that $D(L) < \kappa/2$, so that if $D(H) \leq 1 - \kappa$, then $D(M) \geq \kappa/2$. Consider further partitioning the set M of “medium weight” points into buckets M_1, \dots, M_r where $r = \log_{1+\kappa}(2/\kappa^2) = \Theta(\log(1/\kappa)/\kappa)$ and the bucket M_j is defined as follows:

$$M_j \stackrel{\text{def}}{=} \left\{ i : (1 + \kappa)^{j-1} \cdot \frac{\kappa}{2N} \leq D(i) < (1 + \kappa)^j \cdot \frac{\kappa}{2N} \right\}. \quad (73)$$

We consider the following “desirable” events.

1. Let E_1 be the event that conditioned on the existence of a bucket M_j such that $D(M_j) \geq \kappa/2r = \Omega(\kappa^2/\log(1/\kappa))$, there exists a point $x^* \in X$ that belongs to M_j . By the setting of the size of the sample X , the (conditional) event E_1 holds with probability at least $1 - 1/40$.
2. Let E_2 be the event that all calls to ESTIMATE-NEIGHBORHOOD return an output as specified by Lemma 3. By Lemma 3, the setting of the confidence parameter δ in each call and a union bound over all $|X|$ calls, E_2 holds with probability at least $1 - 1/40$.

Algorithm 12: Procedure FIND-REFERENCE

Input: PCOND and SAMP query access to a distribution D and a parameter $\kappa \in (0, 1/4]$

1. Select a sample X of $\Theta(\log(1/\kappa)/\kappa^2)$ points distributed according to D .
 2. For each $x \in X$ do the following:
 - (a) Call ESTIMATE-NEIGHBORHOOD with the parameters κ as in the input to FIND-REFERENCE, $\beta = \kappa^2/(40 \log(1/\kappa))$, $\eta = \kappa$, and $\delta = 1/(40|X|)$. Let $\theta = \kappa\eta\beta\delta/64 = \Theta(\kappa^6/\log^2(1/\kappa))$ (as in FIND-REFERENCE).
 - (b) If ESTIMATE-NEIGHBORHOOD returns a pair $(\hat{w}(x), \alpha(x))$ such that $\hat{w}(x) < \kappa^2/20 \log(1/\kappa)$, then go to Line 2 and continue with next $x \in X$.
 - (c) Select a sample Y_x of size $\Theta(\log^2(1/\kappa)/\kappa^5)$ distributed uniformly.
 - (d) For each $y \in Y_x$ call COMPARE($\{x\}, \{y\}, \theta/4, 4, 1/40|X||Y_x|$), and let the output be denoted $\rho_x(y)$.
 - (e) Let $\hat{\mu}(x)$ be the fraction of occurrences of $y \in Y_x$ such that $\rho_x(y) \in [1/(1 + \alpha + \theta/2), 1 + \alpha + \theta/2]$.
 - (f) Set $\hat{D}(x) = \hat{w}(x)/(\hat{\mu}(x)N)$.
 3. If for some point $x \in X$ we have $\hat{w}(x) \geq \kappa^2/20 \log(1/\kappa)$, $\hat{\mu}(x) \geq \kappa^3/20 \log(1/\kappa)$, and $\kappa/4n \leq \hat{D}(x) \leq 2/(\kappa N)$, then return $(x, \hat{D}(x))$. Otherwise return No-Pair.
-

3. Let E_3 be the event that for each $x \in X$ we have the following.

- (a) If $\frac{|U_{\alpha(x)}(x)|}{N} \geq \frac{\kappa^3}{40 \log(1/\kappa)}$, then $\frac{|Y_x \cap U_{\alpha(x)}(x)|}{|Y_x|} \in [1 - \eta/2, 1 + \eta/2] \frac{|U_{\alpha(x)}(x)|}{N}$;
If $\frac{|U_{\alpha(x)}(x)|}{N} < \frac{\kappa^3}{40 \log(1/\kappa)}$, then $\frac{|Y_x \cap U_{\alpha(x)}(x)|}{|Y_x|} < \frac{\kappa^3}{30 \log(1/\kappa)}$;
- (b) Let $\Delta_{\alpha(x), \theta}(x) \stackrel{\text{def}}{=} U_{\alpha(x)+\theta}(x) \setminus U_{\alpha(x)}(x)$ (where θ is as specified by the algorithm).
If $\frac{|\Delta_{\alpha(x), \theta}(x)|}{N} \geq \frac{\kappa^4}{240 \log(1/\kappa)}$, then $\frac{|Y_x \cap \Delta_{\alpha(x), \theta}(x)|}{|Y_x|} \leq 2 \cdot \frac{|\Delta_{\alpha(x), \theta}(x)|}{N}$;
If $\frac{|\Delta_{\alpha(x), \theta}(x)|}{N} < \frac{\kappa^4}{240 \log(1/\kappa)}$, then $\frac{|Y_x \cap \Delta_{\alpha(x), \theta}(x)|}{|Y_x|} < \frac{\kappa^4}{120 \log(1/\kappa)}$.

By the size of each set Y_x and a union bound over all $x \in X$, the event E_3 holds with probability at least $1 - 1/40$.

4. Let E_4 be the event that all calls to COMPARE return an output as specified by Lemma 2. By Lemma 2, the setting of the confidence parameter δ in each call and a union bound over all (at most) $|X| \cdot |Y|$ calls, E_3 holds with probability at least $1 - 1/40$.

Assuming events E_1 – E_4 all hold (which occurs with probability at least $9/10$) we have the following.

1. By E_2 , for each $x \in X$ such that $\hat{w}(x) \geq \kappa^2/20 \log(1/\kappa)$ (so that x may be selected for the output of the procedure) we have that $D(U_{\alpha(x)}(x)) \geq \kappa^2/40 \log(1/\kappa)$.

The event E_2 also implies that for each $x \in X$ we have that $D(\Delta_{\alpha(x),\theta(x)}) \leq \eta\beta/16 \leq (\eta/16) \cdot D(U_{\alpha(x)}(x))$, so that

$$\frac{|\Delta_{\alpha(x),\theta(x)}|}{N} \leq \frac{\eta(1+\alpha(x))(1+\alpha(x)+\theta)}{16} \cdot \frac{|U_{\alpha(x)}(x)|}{N} \leq \frac{\eta}{6} \cdot \frac{|U_{\alpha(x)}(x)|}{N}. \quad (74)$$

2. Consider any $x \in X$ such that $\hat{w}(x) \geq \kappa^2/20 \log(1/\kappa)$. Let $T_x \stackrel{\text{def}}{=} \{y : \rho_x(y) \in [1/(1+\alpha+\theta/2), (1+\alpha+\theta/2)]\}$, so that $\hat{\mu}(x) = |T_x|/|Y_x|$. By E_4 , for each $y \in Y_x \cap U_{\alpha(x)}(x)$ we have that $\rho_x(y) \leq (1+\alpha)(1+\theta/4) \leq (1+\alpha+\theta/2)$ and $\rho_x(y) \geq (1+\alpha)^{-1}(1-\theta/4) \geq (1+\alpha+\theta/2)^{-1}$, so that $y \in T_x$. On the other hand, for each $y \notin Y_x \cap U_{\alpha(x)+\theta}(x)$ we have that $\rho_x(y) > (1+\alpha+\theta)(1-\theta/4) \geq 1+\alpha+\theta/2$ or $\rho_x(y) < (1+\alpha+\theta)^{-1}(1-\theta/4) < (1+\alpha+\theta/2)^{-1}$, so that $y \notin T_x$. It follows that

$$Y_x \cap U_{\alpha(x)}(x) \subseteq T_x \subseteq Y_x \cap (U_{\alpha(x)}(x) \cup \Delta_{\alpha(x),\theta(x)}). \quad (75)$$

By E_3 , when $\hat{\mu}(x) = |T_x|/|Y_x| \geq \kappa^3/20 \log(1/\kappa)$, then necessarily $\hat{\mu}(x) \in [1-\eta, 1+\eta]|U_{\alpha(x)}(x)|/N$. To verify this consider the following cases.

- (a) If $\frac{|U_{\alpha(x)}(x)|}{N} \geq \frac{\kappa^3}{40 \log(1/\kappa)}$, then (by the left-hand-side of Equation (75) and the definition of E_3) we get that $\hat{\mu}(x) \geq (1-\eta/2) \frac{|U_{\alpha(x)}(x)|}{N}$, and (by the right-hand-side of Equation (75), Equation (74) and E_3) we get that $\hat{\mu}(x) \leq (1+\eta/2) \frac{|U_{\alpha(x)}(x)|}{N} + 2(\eta/6) \frac{|U_{\alpha(x)}(x)|}{N} < (1+\eta) \frac{|U_{\alpha(x)}(x)|}{N}$.
- (b) If $\frac{|U_{\alpha(x)}(x)|}{N} < \frac{\kappa^3}{40 \log(1/\kappa)}$, then (by the right-hand-side of Equation (75), Equation (74) and E_3) we get that $\hat{\mu}(x) < \frac{\kappa^3}{30 \log(1/\kappa)} + \frac{\kappa^4}{120 \log(1/\kappa)} < \kappa^3/20 \log(1/\kappa)$.
3. If $D(H) \leq 1-\kappa$, so that $D(M) \geq \kappa/2$, then there exists at least one bucket M_j such that $D(M_j) \geq \kappa/2r = \Omega(\kappa^2/\log(1/\kappa))$. By E_1 , the sample X contains a point $x^* \in M_j$. By the definition of the buckets, for this point x^* we have that $D(U_\kappa(x^*)) \geq \kappa/2r \geq \kappa^2/(10 \log(1/\kappa))$ and $|U_\kappa(x^*)| \geq (\kappa^2/2r)N \geq \kappa^3/(10 \log(1/\kappa))$.

By the first two items above and the setting $\eta = \kappa$ we have that for each x such that $\hat{w}(x) \geq \kappa^2/20 \log(1/\kappa)$ and $\hat{\mu}(x) \geq \kappa^3/20 \log(1/\kappa)$,

$$\widehat{D}(x) \in \left[\frac{1-\kappa}{1+\kappa}, \frac{1+\kappa}{1-\kappa} \right] D(x) \subset [1-2\kappa, 1+3\kappa] D(x).$$

Thus, if the algorithm outputs a pair $(x, \widehat{D}(x))$ then it satisfies the condition stated in both items of the lemma. This establishes the second item in the lemma. By combining all three items we get that if $D(H) \geq 1-\kappa$ then the algorithm outputs a pair $(x, \widehat{D}(x))$ (where possibly, but not necessarily, $x = x^*$), and the first item is established as well.

Turning to the query complexity, the total number of PCOND queries performed in the $|X| = O(\log(1/\kappa)/\kappa^2)$ calls to ESTIMATE-NEIGHBORHOOD is $O\left(\frac{|X| \log(1/\gamma)^2 \log(1/(\beta\eta))}{\kappa^2 \eta^4 (\beta)^3 \gamma^2}\right) = \tilde{O}(1/\kappa^{18})$, and the total number of PCOND queries performed in the calls to COMPARE (for at most all pairs $x \in X$ and $y \in Y_x$) is $\tilde{O}(1/\kappa^{20})$. ■

8 Algorithms and lower bounds for testing uniformity with ICOND_D

In this section we consider ICOND algorithms for testing whether an unknown distribution D over $[N]$ is the uniform distribution versus ϵ -far from uniform. Our results show that ICOND algorithms are not as powerful as PCOND algorithms for this basic testing problem; we give a $\text{poly}(\log N, 1/\epsilon)$ -query ICOND_D algorithm, and prove that any ICOND_D algorithm must make $\tilde{\Omega}(\log N)$ queries.

8.1 A $\tilde{O}\left(\frac{\log^3 N}{\epsilon^3}\right)$ -query ICOND_D algorithm for testing uniformity

In this subsection we describe an algorithm $\text{ICOND}_D\text{-TEST-UNIFORM}$ and prove the following theorem:

Theorem 15 $\text{ICOND}_D\text{-TEST-UNIFORM}$ is a $\tilde{O}\left(\frac{\log^3 N}{\epsilon^3}\right)$ -query ICOND_D testing algorithm for uniformity, i.e. it outputs ACCEPT with probability at least $2/3$ if $D = \mathcal{U}$ and outputs REJECT with probability at least $2/3$ if $d_{\text{TV}}(D, \mathcal{U}) \geq \epsilon$.

Intuition. Recall that, as mentioned in Section 4.1, any distribution D which is ϵ -far from uniform must put $\Omega(\epsilon)$ probability mass on “significantly heavy” elements (that is, if we define $H' = \{h \in [N] \mid D(h) \geq \frac{1}{N} + \frac{\epsilon}{4N}\}$, it must hold that $D(H') \geq \epsilon/4$). Consequently a sample of $O(1/\epsilon)$ points drawn from D will contain such a point with high probability. Thus, a natural approach to testing whether D is uniform is to devise a procedure that, given an input point y , can distinguish between the case that $y \in H'$ and the case that $D(y) = 1/N$ (as it is when $D = \mathcal{U}$).

We give such a procedure, which uses the ICOND_D oracle to perform a sort of binary search over intervals. The procedure successively “weighs” narrower and narrower intervals until it converges on the single point y . In more detail, we consider the *interval tree* whose root is the whole domain $[N]$, with two children $\{1, \dots, N/2\}$ and $\{N/2 + 1, \dots, N\}$, and so on, with a single point at each of the N leaves. Our algorithm starts at the root of the tree and goes down the path that corresponds to y ; at each child node it COMPARE to compare the weight of the current node to the weight of its sibling under D . If at any point the estimate deviates significantly from the value it should have if D were uniform (namely the weights should be essentially equal, with slight deviations because of even/odd issues), then the algorithm rejects. Assuming the algorithm does not reject, it provides a $(1 \pm O(\epsilon))$ -accurate multiplicative estimate of $D(y)$, and the algorithm checks whether this estimate is sufficiently close to $1/N$ (rejecting if this is not the case). If no point in a sample of $\Theta(1/\epsilon)$ points (drawn according to D) causes rejection, then the algorithm accepts.

The algorithm we use to perform the “binary search” described above is Algorithm 13, BINARY-DESCENT . We begin by proving correctness for it:

Lemma 27 Suppose the algorithm BINARY-DESCENT is run with inputs $\epsilon \in (0, 1]$, $a = 1$, $b = N$, and $y \in [N]$, and is provided ICOND oracle access to distribution D over $[N]$. It performs $\tilde{O}(\log^3 N/\epsilon^2)$ queries and either outputs a value $\hat{D}(y)$ or REJECT, where the following holds:

Algorithm 13: BINARY-DESCENT

Input: parameter $\epsilon > 0$; integers $1 \leq a \leq b \leq N$; $y \in [a, b]$; query access to ICOND_D oracle

- 1: **if** $a = b$ **then**
- 2: **return** 1
- 3: **end if**
- 4: Let $c = \lfloor \frac{a+b}{2} \rfloor$; $\Delta = (b - a + 1)/2$.
- 5: **if** $y \leq c$ **then**
- 6: Define $I_y = \{a, \dots, c\}$, $I_{\bar{y}} = \{c + 1, \dots, b\}$ and $\rho = \lceil \Delta \rceil / \lfloor \Delta \rfloor$
- 7: **else**
- 8: Define $I_{\bar{y}} = \{a, \dots, c\}$, $I_y = \{c + 1, \dots, b\}$ and $\rho = \lfloor \Delta \rfloor / \lceil \Delta \rceil$
- 9: **end if**
- 10: Call COMPARE on $I_y, I_{\bar{y}}$ with parameters $\eta = \frac{\epsilon}{48 \log N}$, $K = 2$, $\delta = \frac{\epsilon}{100(1 + \log N)}$ to get an estimate $\hat{\rho}$ of $D(I_y)/D(I_{\bar{y}})$
- 11: **if** $\hat{\rho} \notin [1 - \frac{\epsilon}{48 \log N}, 1 + \frac{\epsilon}{48 \log N}] \cdot \rho$ (this includes the case that $\hat{\rho}$ is High or Low) **then**
- 12: **return** REJECT
- 13: **end if**
- 14: Call recursively BINARY-DESCENT on input $(\epsilon, \text{the endpoints of } I_y, y)$;
- 15: **if** BINARY-DESCENT returns a value ν **then**
- 16: **return** $\frac{\hat{\rho}}{1 + \hat{\rho}} \cdot \nu$
- 17: **else**
- 18: **return** REJECT
- 19: **end if**

Algorithm 14: ICOND_D-TEST-UNIFORM

Input: error parameter $\epsilon > 0$; query access to ICOND_D oracle

- 1: Draw $t = \frac{20}{\epsilon}$ points y_1, \dots, y_t from SAMP_D .
- 2: **for** $j = 1$ to t **do**
- 3: Call BINARY-DESCENT($\epsilon, 1, N, y_j$) and return REJECT if it rejects, otherwise let \hat{d}_j be the value it returns as its estimate of $D(y_j)$
- 4: **if** $\hat{d}_j \notin [1 - \frac{\epsilon}{12}, 1 + \frac{\epsilon}{12}] \cdot \frac{1}{N}$ **then**
- 5: **return** REJECT
- 6: **end if**
- 7: **end for**
- 8: **return** ACCEPT

1. if $D(y) \geq \frac{1}{N} + \frac{\epsilon}{4N}$, then with probability at least $1 - \frac{\epsilon}{100}$ the procedure either outputs a value $\hat{D}(y) \in [1 - \epsilon/12, 1 + \epsilon/12]D(y)$ or REJECT;
2. if $D = \mathcal{U}$, then with probability at least $1 - \frac{\epsilon}{100}$ the procedure outputs a value $\hat{D}(y) \in [1 - \epsilon/12, 1 + \epsilon/12] \cdot \frac{1}{N}$.

Proof of Lemma 27: The claimed query bound is easily verified, since the recursion depth is at most $1 + \log N$ and the only queries made are during calls to COMPARE, each of which performs $O(\log(1/\delta)/\gamma^2) = \tilde{O}(\log^2 N/\epsilon^2)$ queries.

Let E_0 be the event that all calls to COMPARE satisfy the conditions in Lemma 2; since each of them succeeds with probability at least $1 - \delta = 1 - \frac{\epsilon}{100(1+\log N)}$, a union bound shows that E_0 holds with probability at least $1 - \epsilon/100$. We hereafter condition on E_0 .

We first prove the second part of the lemma where $D = \mathcal{U}$. Fix any specific recursive call, say the j -th, during the execution of the procedure. The intervals $I_y^{(j)}, I_{\bar{y}}^{(j)}$ used in that execution of the algorithm are easily seen to satisfy $D(I_y)/D(I_{\bar{y}}) \in [1/K, K]$ (for $K = 2$), so by event E_0 it must be the case that COMPARE returns an estimate $\hat{\rho}_j \in [1 - \frac{\epsilon}{48 \log N}, 1 + \frac{\epsilon}{48 \log N}] \cdot D(I_y^{(j)})/D(I_{\bar{y}}^{(j)})$. Since $D = \mathcal{U}$, we have that $D(I_y^{(j)})/D(I_{\bar{y}}^{(j)}) = \rho^{(j)}$, so the overall procedure returns a numerical value rather than REJECT.

Let $M = \lceil \log N \rceil$ be the number of recursive calls (i.e., the number of executions of Line 14). Note that we can write $D(y)$ as a product

$$D(y) = \prod_{j=1}^M \frac{D(I_y^{(j)})}{D(I_y^{(j)}) + D(I_{\bar{y}}^{(j)})} = \prod_{j=1}^M \frac{D(I_y^{(j)})/D(I_{\bar{y}}^{(j)})}{D(I_y^{(j)})/D(I_{\bar{y}}^{(j)}) + 1}. \quad (76)$$

We next observe that for any $0 \leq \epsilon' < 1$ and $\rho, d > 0$, if $\hat{\rho} \in [1 - \epsilon', 1 + \epsilon']d$ then we have $\frac{\hat{\rho}}{\hat{\rho}+1} \in [1 - \frac{\epsilon'}{2}, 1 + \epsilon']\frac{d}{d+1}$ (by straightforward algebra). Applying this M times, we get

$$\begin{aligned} \prod_{j=1}^M \frac{\hat{\rho}_j}{\hat{\rho}_j + 1} &\in \left[\left(1 - \frac{\epsilon}{96 \log N}\right)^M, \left(1 + \frac{\epsilon}{48 \log N}\right)^M \right] \cdot \prod_{j=1}^M \frac{D(I_y^{(j)})/D(I_{\bar{y}}^{(j)})}{D(I_y^{(j)})/D(I_{\bar{y}}^{(j)}) + 1} \\ &\in \left[\left(1 - \frac{\epsilon}{96 \log N}\right)^M, \left(1 + \frac{\epsilon}{48 \log N}\right)^M \right] \cdot D(y) \\ &\in \left[1 - \frac{\epsilon}{12}, 1 + \frac{\epsilon}{12} \right] D(y). \end{aligned}$$

Since $\prod_{j=1}^M \frac{\hat{\rho}_j}{\hat{\rho}_j + 1}$ is the value that the procedure outputs, the second part of the lemma is proved.

The proof of the first part of the lemma is virtually identical. The only difference is that now it is possible that COMPARE outputs High or Low at some call (since D is not uniform it need not be the case that $D(I_y^{(j)})/D(I_{\bar{y}}^{(j)}) = \rho^{(j)}$), but this is not a problem for (i) since in that case BINARY-DESCENT would output REJECT. ■

See Algorithm 13 for a description of the testing algorithm ICOND_D-TEST-UNIFORM. We now prove Theorem 15:

Proof of Theorem 15: Define E_1 to be the event that all calls to BINARY-DESCENT satisfy the conclusions of Lemma 27. With a union bound over all these $t = 20/\epsilon$ calls, we have $\Pr[E_1] \geq 8/10$.

Completeness: Suppose $D = \mathcal{U}$, and condition again on E_1 . Since this implies that BINARY-DESCENT will always return a value, the only case ICOND_D-TEST-UNIFORM might reject is by reaching Line 5. However, since it is the case that every value \hat{d}_j returned by the procedure satisfies $\hat{D}(y) \in [1 - \epsilon/12, 1 + \epsilon/12] \cdot \frac{1}{N}$, this can never happen.

Soundness: Suppose $d_{\text{TV}}(D, \mathcal{U}) \geq \epsilon$. Let E_2 be the event that at least one of the y_i 's drawn in Line 1 belongs to H' . As $D(H') \geq \epsilon/4$, we have $\Pr[E_2] \geq 1 - (1 - \epsilon/4)^{20/\epsilon} \geq 9/10$. Conditioning on both E_1 and E_2 , for such a y_j , one of two cases below holds:

- either the call to BINARY-DESCENT outputs REJECT and ICOND_D-TEST-UNIFORM outputs REJECT;
- or a value \hat{d}_j is returned, for which $\hat{d}_j \geq (1 - \frac{\epsilon}{12})(1 + \frac{\epsilon}{4}) \cdot \frac{1}{N} > (1 + \epsilon/12)/N$ (where we used the fact that E_1 holds); and ICOND_D-TEST-UNIFORM reaches Line 5 and rejects.

Since $\Pr[E_1 \cup E_2] \geq 7/10$, ICOND_D-TEST-UNIFORM is correct with probability at least 2/3. Finally, the claimed query complexity directly follows from the $t = \Theta(1/\epsilon)$ calls to BINARY-DESCENT, each of which makes $\tilde{O}(\log^3 N/\epsilon^2)$ queries to ICOND_D. ■

8.2 An $\Omega(\log N/\log \log N)$ lower bound for ICOND_D algorithms that test uniformity

This subsection proves that any ICOND_D algorithm that ϵ -tests uniformity even for constant ϵ must have query complexity $\tilde{\Omega}(\log N)$. This shows that our algorithm in the previous subsection is not too far from optimal, and sheds light on a key difference between ICOND and PCOND oracles.

Theorem 16 *Fix $\epsilon = 1/3$. Any ICOND_D algorithm for testing whether $D = \mathcal{U}$ versus $d_{\text{TV}}(D, D^*) \geq \epsilon$ must make $\Omega\left(\frac{\log N}{\log \log N}\right)$ queries.*

To prove this lower bound we define a probability distribution \mathcal{P}_{No} over possible “No”-distributions (i.e. distributions that have variation distance at least 1/3 from \mathcal{U}). A distribution drawn from \mathcal{P}_{No} is constructed as follows: first, we partition $[N]$ into $b = 2^X$ consecutive intervals of the same size $\Delta = \frac{N}{2^X}$, which we refer to as “blocks”, where X is a random variable distributed uniformly on the set $\{\frac{1}{3} \log N, \frac{1}{3} \log N + 1, \dots, \frac{2}{3} \log N\}$. Once the block size Δ is determined, a random offset y is drawn uniformly at random in $[N]$, and all block endpoints are shifted by y modulo $[N]$ (intuitively, this prevents the testing algorithm from “knowing” a priori that specific points are endpoints of blocks). Finally, independently for each block, a fair coin is thrown to determine its *profile*: with probability 1/2, each point in the first half of the block will have probability weight $\frac{1-2\epsilon}{N}$ and each point in the second half will have probability $\frac{1+2\epsilon}{N}$ (such a block is said to be a *low-high* block, with profile $\downarrow\uparrow$). With probability 1/2 the reverse is true: each point in the first half has probability $\frac{1+2\epsilon}{N}$ and each point in the second half has probability $\frac{1-2\epsilon}{N}$ (a *high-low* block $\uparrow\downarrow$). It is clear that each distribution D in the support of \mathcal{P}_{No} defined in this way indeed has $d_{\text{TV}}(D, \mathcal{U}) = \epsilon$.

To summarize, each “No”-distribution D in the support of \mathcal{P}_{No} is parameterized by $(b + 2)$ parameters: its block size Δ , offset y , and profile $\vartheta \in \{\downarrow\uparrow, \uparrow\downarrow\}^b$. Note that regardless of the profile vector, each block always has weight exactly Δ/N .

We note that while there is only one “Yes”-distribution \mathcal{U} , it will sometimes be convenient for the analysis to think of \mathcal{U} as resulting from the same initial process of picking a block size and offset, but without the subsequent choice of a profile vector. We sometimes refer to this as the “fake construction” of the uniform distribution \mathcal{U} (the reason for this will be clear later).

The proof of Theorem 16 will be carried out in two steps. First we shall restrict the analysis to *non-adaptive algorithms*, and prove the lower bound for such algorithms. This result will then be extended to the general setting by introducing (similarly to Section 5.2) the notion of a *query-faking algorithm*, and reducing the behavior of adaptive algorithms to non-adaptive ones through an appropriate sequence of such query-faking algorithms.

Before proceeding, we define the *transcript* of the interaction between an algorithm and a ICOND_D oracle. Informally, the transcript captures the entire history of interaction between the algorithm and the ICOND_D oracle during the whole sequence of queries.

Definition 7 Fix any (possibly adaptive) testing algorithm A that queries an ICOND_D oracle. The transcript of A is a sequence $\mathcal{T} = (I_i, s_i)_{i \in \mathbb{N}^*}$ of pairs, where I_i is the i -th interval provided by the algorithm as input to ICOND_D , and $s_i \in I_i$ is the response that ICOND_D provides to this query. Given a transcript \mathcal{T} , we shall denote by $\mathcal{T}|_k$ the partial transcript induced by the first k queries, i.e. $\mathcal{T}|_k = (I_i, s_i)_{1 \leq i \leq k}$.

Equipped with these definitions, we now turn to proving the theorem in the special case of non-adaptive testing algorithms. Observe that there are three different sources of randomness in our arguments: (i) the draw of the “No”-instance from \mathcal{P}_{No} , (ii) the internal randomness of the testing algorithm; and (iii) the random draws from the oracle. Whenever there could be confusion we shall explicitly state which probability space is under discussion.

8.2.1 Against non-adaptive algorithms

Throughout this subsection we assume that A is an arbitrary, fixed, non-adaptive, randomized algorithm that makes exactly $q \leq \tau \cdot \frac{\log N}{\log \log N}$ queries to ICOND_D ; here $\tau > 0$ is some absolute constant that will be determined in the course of the analysis. (The assumption that A always makes exactly q queries is without loss of generality since if in some execution the algorithm makes $q' < q$ queries, it can perform additional “dummy” queries). In this setting algorithm A corresponds to a distribution P_A over q -tuples $\bar{I} = (I_1, \dots, I_q)$ of query intervals. The following theorem will directly imply Theorem 16 in the case of non-adaptive algorithms:

Theorem 17

$$\left| \Pr_{D \sim \mathcal{P}_{\text{No}}} [A^{\text{ICOND}_D} \text{ outputs ACCEPT}] - \Pr [A^{\text{ICOND}_u} \text{ outputs ACCEPT}] \right| \leq 1/5. \quad (77)$$

Observe that in the first probability of Equation (77) the randomness is taken over the draw of D from \mathcal{P}_{No} , the draw of $\bar{I} \sim P_A$ that A performs to select its sequence of query intervals, and the randomness of the ICOND_D oracle. In the second one the randomness is just over the draw of \bar{I} from P_A and the randomness of the $\text{ICOND}_{\mathcal{U}}$ oracle.

Intuition for Theorem 17. The high-level idea is that the algorithm will not be able to distinguish between the uniform distribution and a “No”-distribution unless it manages to learn something about the “structure” of the blocks in the “No”-case, either by guessing (roughly) the right block size, or by guessing (roughly) the location of a block endpoint and querying a short interval containing such an endpoint.

In more detail, we define the following “bad events” (over the choice of D and the points s_i) for a fixed sequence $\bar{I} = (I_1, \dots, I_q)$ of queries (the dependence on \bar{I} is omitted in the notation for the sake of readability):

$$\begin{aligned} B_{\text{size}}^{\text{N}} &= \{ \exists i \in [q] \mid \Delta / \log N \leq |I_i| \leq \Delta \cdot (\log N)^2 \} \\ B_{\text{boundary}}^{\text{N}} &= \{ \exists i \in [q] \mid |I_i| < \Delta / \log N \text{ and } I_i \text{ intersects two blocks} \} \\ B_{i,\text{outer}}^{\text{N}} &= \{ \Delta \cdot (\log N)^2 < |I_i| \text{ and } s_i \text{ belongs to a block not contained entirely in } I_i \} \quad i \in [q] \\ B_{i,\text{collide}}^{\text{N}} &= \{ \Delta \cdot (\log N)^2 < |I_i| \text{ and } \exists j < i, s_i \text{ and } s_j \text{ belong to the same block} \} \quad i \in [q] \end{aligned}$$

The first two events depend only on the draw of D from \mathcal{P}_{No} which determines Δ and y , while the last $2q$ events also depend on the random draws of s_i from the ICOND_D oracle. We define in the same fashion the corresponding bad events for the “Yes”-instance (i.e. the uniform distribution \mathcal{U}) $B_{\text{size}}^{\text{Y}}$, $B_{\text{boundary}}^{\text{Y}}$, $B_{i,\text{outer}}^{\text{Y}}$ and $B_{i,\text{collide}}^{\text{Y}}$, using the notion of the “fake construction” of \mathcal{U} mentioned above.

Events $B_{\text{size}}^{\text{N}}$ and $B_{\text{size}}^{\text{Y}}$ correspond to the possibility, mentioned above, that algorithm A “guesses” essentially the right block size, and events $B_{\text{boundary}}^{\text{N}}$ and $B_{\text{boundary}}^{\text{Y}}$ correspond to the possibility that algorithm A “guesses” a short interval containing a block endpoint. The final bad events correspond to A guessing a “too-large” block size but “getting lucky” with the sample returned by ICOND , either because the sample belongs to one of the (at most two) outer blocks not entirely contained in the query interval, or because A has already received a sample from the same block as the current sample.

We can now describe the *failure events* for both the uniform distribution and for a “No”-distribution as the union of the corresponding bad events:

$$\begin{aligned} B_{(\bar{I})}^{\text{N}} &= B_{\text{size}}^{\text{N}} \cup B_{\text{boundary}}^{\text{N}} \cup \left(\bigcup_{i=1}^q B_{i,\text{outer}}^{\text{N}} \right) \cup \left(\bigcup_{i=1}^q B_{i,\text{collide}}^{\text{N}} \right) \\ B_{(\bar{I})}^{\text{Y}} &= B_{\text{size}}^{\text{Y}} \cup B_{\text{boundary}}^{\text{Y}} \cup \left(\bigcup_{i=1}^q B_{i,\text{outer}}^{\text{Y}} \right) \cup \left(\bigcup_{i=1}^q B_{i,\text{collide}}^{\text{Y}} \right) \end{aligned}$$

These failure events can be interpreted, from the point of view of the algorithm A , as the “opportunity to potentially learn something;” we shall argue below that if the failure events do not occur then the algorithm gains no information about whether it is interacting with the uniform distribution or with a “No”-distribution.

Structure of the proof of Theorem 17. First, observe that since the transcript is the result of the interaction of the algorithm and the oracle on a randomly chosen distribution, it is itself a random variable; we will be interested in the distribution over this random variable induced by the draws from the oracle and the choice of D . More precisely, for a fixed sequence of query sets \bar{I} , let $Z_{\bar{I}}^N$ denote the random variable over “No”-transcripts generated when D is drawn from \mathcal{P}_{No} . Note that this is a random variable over the probability space defined by the random draw of D and the draws of s_i by $\text{ICOND}_D(I_i)$. We define $\mathfrak{A}_{\bar{I}}^N$ as the resulting distribution over these “No”-transcripts. Similarly, $Z_{\bar{I}}^Y$ will be the random variable over “Yes”-transcripts, with corresponding distribution $\mathfrak{A}_{\bar{I}}^Y$.

As noted earlier, the nonadaptive algorithm A corresponds to a distribution P_A over q -tuples \bar{I} of query intervals. We define \mathfrak{A}^N as the distribution over transcripts corresponding to first drawing \bar{I} from P_A and then making a draw from $\mathfrak{A}_{\bar{I}}^N$. Similarly, we define \mathfrak{A}^Y as the distribution over transcripts corresponding to first drawing \bar{I} from P_A and then making a draw from $\mathfrak{A}_{\bar{I}}^Y$.

To prove Theorem 17 it is sufficient to show that the two distributions over transcripts described above are statistically close:

Lemma 28 $d_{\text{TV}}(\mathfrak{A}^Y, \mathfrak{A}^N) \leq 1/5$.

The proof of this lemma is structured as follows: first, for any *fixed* sequence of q queries \bar{I} , we bound the probability of the failure events, both for the uniform and the “No”-distributions:

Claim 29 *For each fixed sequence \bar{I} of q query intervals, we have*

$$\Pr[B_{(\bar{I})}^Y] \leq 1/10 \quad \text{and} \quad \Pr_{D \leftarrow \mathcal{P}_{\text{No}}}[B_{(\bar{I})}^N] \leq 1/10.$$

(Note that the first probability above is taken only over the randomness of the ICOND_U responses, while the second is over the random draw of $D \sim \mathcal{P}_{\text{No}}$ and over the ICOND_D responses.)

Next we show that, provided the failure events do not occur, the distribution over transcripts is exactly the same in both cases:

Claim 30 *Fix any sequence $\bar{I} = (I_1, \dots, I_q)$ of q queries. Then, conditioned on their respective failure events not happening, $Z_{\bar{I}}^N$ and $Z_{\bar{I}}^Y$ are identically distributed:*

$$\text{for every transcript } \mathcal{T} = ((I_1, s_1), \dots, (I_q, s_q)), \quad \Pr[Z_{\bar{I}}^N = \mathcal{T} \mid \overline{B_{(\bar{I})}^N}] = \Pr[Z_{\bar{I}}^Y = \mathcal{T} \mid \overline{B_{(\bar{I})}^Y}].$$

Finally we combine these two claims to show that the two overall distributions of transcripts are statistically close:

Claim 31 *Fix any sequence of q queries $\bar{I} = (I_1, \dots, I_q)$. Then $d_{\text{TV}}(\mathfrak{A}_{\bar{I}}^N, \mathfrak{A}_{\bar{I}}^Y) \leq 1/5$.*

Lemma 28 (and thus Theorem 17) directly follows from Claim 31 since, using the notation $\bar{s} = (s_1, \dots, s_q)$ for a sequence of q answers to a sequence $\bar{I} = (I_1, \dots, I_q)$ of q queries, which together define a transcript $\mathcal{T}(\bar{I}, \bar{s}) = ((I_1, s_1), \dots, (I_q, s_q))$,

$$\begin{aligned}
d_{\text{TV}}(\mathfrak{A}^Y, \mathfrak{A}^N) &= \frac{1}{2} \sum_{\bar{I}} \sum_{\bar{s}} |P_A(\bar{I}) \cdot \Pr[Z_{\bar{I}}^Y = \mathcal{T}(\bar{I}, \bar{s})] - P_A(\bar{I}) \cdot \Pr[Z_{\bar{I}}^N = \mathcal{T}(\bar{I}, \bar{s})]| \\
&= \sum_{\bar{I}} \sum_{\bar{s}} P_A(\bar{I}) \cdot \sum_{\bar{s}} |\Pr[Z_{\bar{I}}^Y = \mathcal{T}(\bar{I}, \bar{s})] - \Pr[Z_{\bar{I}}^N = \mathcal{T}(\bar{I}, \bar{s})]| \\
&\leq \max_{\bar{I}} \{d_{\text{TV}}(\mathfrak{A}_{\bar{I}}^Y, \mathfrak{A}_{\bar{I}}^N)\} \leq 1/5.
\end{aligned} \tag{78}$$

This concludes the proof of Lemma 28 modulo the proofs of the above claims; we give those proofs in Section 8.2.2 below.

8.2.2 Proof of Claims 29 to 31

To prove Claim 29 we bound the probability of each of the bad events separately, starting with the “No”-case.

- (i) Defining the event $B_{i,\text{size}}^N$ as

$$B_{i,\text{size}}^N = \{\Delta/\log N \leq |I_i| \leq \Delta \cdot (\log N)^2\},$$

we can use a union bound to get $\Pr[B_{\text{size}}^N] \leq \sum_{i=1}^q \Pr[B_{i,\text{size}}^N]$. For any fixed setting of I_i there are $O(\log \log N)$ values of $\Delta \in \{\frac{N}{2^X} \mid X \in \{\frac{1}{3} \log N, \dots, \frac{2}{3} \log N\}\}$ for which $\Delta/\log N \leq |I_i| \leq \Delta \cdot (\log N)^2$. Hence we have $\Pr[B_{i,\text{size}}^N] = O((\log \log N)/\log N)$, and consequently $\Pr[B_{\text{size}}^N] = O(q(\log \log N)/\log N)$.

- (ii) Similarly, defining the event $B_{i,\text{boundary}}^N$ as

$$B_{i,\text{boundary}}^N = \{|I_i| < \Delta/\log N \text{ and } I_i \text{ intersects two blocks}\},$$

we have $\Pr[B_{\text{boundary}}^N] \leq \sum_{i=1}^q \Pr[B_{i,\text{boundary}}^N]$. For any fixed setting of I_i , recalling the choice of a uniform random offset $y \in [N]$ for the blocks, we have that $\Pr[B_{i,\text{boundary}}^N] \leq O(1/\log N)$, and consequently $\Pr[B_{\text{boundary}}^N] = O(q/\log N)$.

- (iii) Fix $i \in [q]$ and recall that $B_{i,\text{outer}}^N = \{\Delta \cdot (\log N)^2 < |I_i| \text{ and } s_i \text{ is drawn from a block } \subsetneq I_i\}$. Fix any outcome for Δ such that $\Delta \cdot (\log N)^2 < |I_i|$ and let us consider only the randomness over the draw of s_i from I_i . Since there are $\Omega((\log N)^2)$ blocks contained entirely in I_i , the probability that s_i is drawn from a block not contained entirely in I_i (there are at most two such blocks, one at each end of I_i) is $O(1/(\log N)^2)$. Hence we have $\Pr[B_{i,\text{outer}}^N] \leq O(1)/(\log N)^2$.

- (iv) Finally, recall that

$$B_{i,\text{collide}}^N = \{\Delta \cdot (\log N)^2 < |I_i| \text{ and } \exists j < i \text{ s.t. } s_i \text{ and } s_j \text{ belong to the same block}\}.$$

Fix $i \in [q]$ and a query interval I_i . Let r_i be the number of blocks in I_i within which resides some previously sampled point s_j , $j \in [i-1]$. Since there are $\Omega((\log N)^2)$ blocks in I_i and $r_i \leq i-1$, the probability that s_i is drawn from a block containing any s_j , $j < i$, is $O(i/(\log N)^2)$. Hence we have $\Pr[B_{i,\text{collide}}^N] = O(i/(\log N)^2)$.

With these probability bounds for bad events in hand, we can prove Claim 29:

Proof of Claim 29: Recall that $q \leq \tau \cdot \frac{\log N}{\log \log N}$. Recalling the definition of $B_{\bar{I}}^N$, a union bound yields

$$\begin{aligned} \Pr[B_{\bar{I}}^N] &\leq \Pr[B_{\text{size}}^N] + \Pr[B_{\text{boundary}}^N] + \sum_{i=1}^q \Pr[B_{i,\text{outer}}^N] + \sum_{i=1}^q \Pr[B_{i,\text{collide}}^N] \\ &= O\left(\frac{q \cdot \log \log N}{\log N}\right) + O\left(\frac{q}{\log N}\right) + \sum_{i=1}^q O\left(\frac{1}{(\log N)^2}\right) + \sum_{i=1}^q O\left(\frac{i}{(\log N)^2}\right) \\ &\leq \frac{1}{10}, \end{aligned}$$

where the last inequality holds for a sufficiently small choice of the absolute constant τ .

The same analysis applies unchanged for $\Pr[B_{\text{size}}^Y]$, and $\Pr[B_{\text{boundary}}^Y]$, using the “fake construction” view of \mathcal{U} as described earlier. The arguments for $\Pr[B_{i,\text{outer}}^Y]$ and $\Pr[B_{i,\text{collide}}^Y]$ go through unchanged as well, and Claim 29 is proved. \blacksquare

Proof of Claim 30: Fix any $\bar{I} = (I_1, \dots, I_q)$ and any transcript $\mathcal{T} = ((I_1, s_1), \dots, (I_q, s_q))$. Recall that the length- ℓ partial transcript $\mathcal{T}|_\ell$ is defined to be $((I_1, s_1), \dots, (I_\ell, s_\ell))$. We define the random variables $Z_{\bar{I},\ell}^N$ and $Z_{\bar{I},\ell}^Y$ to be the length- ℓ prefixes of $Z_{\bar{I}}^N$ and $Z_{\bar{I}}^Y$ respectively. We prove Claim 30 by establishing the following, which we prove by induction on ℓ :

$$\Pr\left[Z_{\bar{I}}^N = \mathcal{T} \mid \overline{B_{\bar{I}}^N}\right] = \Pr\left[Z_{\bar{I}}^Y = \mathcal{T} \mid \overline{B_{\bar{I}}^Y}\right]. \quad (79)$$

For the base case, it is clear that (79) holds with $\ell = 0$. For the inductive step, suppose (79) holds for all $k \in [\ell-1]$. When querying I_ℓ at the ℓ -th step, one of the following cases must hold (since we conditioned on the “bad events” not happening):

- (1) I_ℓ is contained within a half-block (more precisely, either entirely within the first half of a block or entirely within the second half). In this case the “yes” and “no” distribution oracles behave exactly the same since both generate s_i by sampling uniformly from I_i .
- (2) I_ℓ contains many blocks and s_ℓ belongs to a block, contained entirely in I_i , which is “fresh” in the sense that it contains no s_j , $j < i$. In the “No”-case this block may either be high-low or low-high; but since both outcomes have the same probability, there is another transcript with equal probability in which the two profiles are switched. Consequently (over the randomness in the draw of $D \sim \mathcal{P}_{\text{No}}$) the probability of picking s_ℓ in the “No”-distribution case is the same as in the uniform distribution case (i.e, uniform on I_ℓ).
- (3) I_ℓ is contained within one block, but not within one half-block (i.e. I_ℓ intersects both the first and second halves of its block). By the same symmetry argument as in (2), the profile of this block could have been switched, and hence the distribution of s_i is uniform over I_i .

This concludes the proof of Claim 30. \blacksquare

Proof of Claim 31: Given Claims 29 and 30, Claim 31 is an immediate consequence of the following basic fact:

Fact 32 *Let D_1, D_2 be two distributions over the same finite set X . Let E_1, E_2 , be two events such that $D_i[E_i] = \alpha_i \leq \alpha$ for $i = 1, 2$ and the conditional distributions $(D_i)_{\overline{E_i}}$ are identical, i.e. $d_{\text{TV}}((D_1)_{\overline{E_1}}, (D_2)_{\overline{E_2}}) = 0$. Then $d_{\text{TV}}(D_1, D_2) \leq \alpha$.*

Proof: We first observe that since $(D_2)_{\overline{E_2}}(E_2) = 0$ and $(D_1)_{\overline{E_1}}$ is identical to $(D_2)_{\overline{E_2}}$, it must be the case that $(D_1)_{\overline{E_1}}(E_2) = 0$, and likewise $(D_2)_{\overline{E_2}}(E_1) = 0$. This implies that $D_1(E_2 \setminus E_1) = D_2(E_1 \setminus E_2) = 0$. Now let us write

$$\begin{aligned} 2d_{\text{TV}}(D_1, D_2) &= \sum_{x \in X \setminus (E_1 \cup E_2)} |D_1(x) - D_2(x)| + \sum_{x \in E_1 \cap E_2} |D_1(x) - D_2(x)| + \\ &\quad \sum_{x \in E_1 \setminus E_2} |D_1(x) - D_2(x)| + \sum_{x \in E_2 \setminus E_1} |D_1(x) - D_2(x)|. \end{aligned}$$

We may upper bound $\sum_{x \in E_1 \cap E_2} |D_1(x) - D_2(x)|$ by $\sum_{x \in E_1 \cap E_2} (D_1(x) + D_2(x)) = D_1(E_1 \cap E_2) + D_2(E_1 \cap E_2)$, and the above discussion gives $\sum_{x \in E_1 \setminus E_2} |D_1(x) - D_2(x)| = D_1(E_1 \setminus E_2)$ and $\sum_{x \in E_2 \setminus E_1} |D_1(x) - D_2(x)| = D_2(E_2 \setminus E_1)$. We thus have

$$\begin{aligned} 2d_{\text{TV}}(D_1, D_2) &\leq \sum_{x \in X \setminus (E_1 \cup E_2)} |D_1(x) - D_2(x)| + D_1(E_1) + D_2(E_2) \\ &\leq \sum_{x \in X \setminus (E_1 \cup E_2)} |D_1(x) - D_2(x)| + \alpha_1 + \alpha_2. \end{aligned}$$

Finally, since $d_{\text{TV}}((D_1)_{\overline{E_1}}, (D_2)_{\overline{E_2}}) = 0$, we have

$$\begin{aligned} \sum_{x \in X \setminus (E_1 \cup E_2)} |D_1(x) - D_2(x)| &= |D_1(X \setminus (E_1 \cup E_2)) - D_2(X \setminus (E_1 \cup E_2))| \\ &= |D_1(\overline{E_1}) - D_2(\overline{E_2})| = |\alpha_1 - \alpha_2|. \end{aligned}$$

Thus $2d_{\text{TV}}(D_1, D_2) \leq |\alpha_1 - \alpha_2| + \alpha_1 + \alpha_2 = 2 \max\{\alpha_1, \alpha_2\} \leq 2\alpha$, and the fact is established. \blacksquare

This concludes the proof of Claim 31. \blacksquare

8.2.3 A lower bound against adaptive algorithms: Proof of Theorem 16

Throughout this subsection A denotes a general adaptive algorithm that makes $q \leq \tau \cdot \frac{\log N}{\log \log N}$ queries, where as before $\tau > 0$ is an absolute constant. Theorem 16 is a consequence of the following theorem, which deals with adaptive algorithms:

Theorem 18

$$\left| \Pr_{D \leftarrow \mathcal{P}_{\text{No}}} [A^{\text{ICOND}_D} \text{ outputs ACCEPT}] - \Pr [A^{\text{ICOND}_u} \text{ outputs ACCEPT}] \right| \leq 1/4. \quad (80)$$

The idea here is to use the previous analysis for non-adaptive algorithms, and argue that “adaptiveness does not really help” to distinguish between $D = \mathcal{U}$ and $D \sim \mathcal{P}_{\text{No}}$ given access to ICOND_D .

As in Section 5.2, we will need the notion of an algorithm *faking queries*. Given an adaptive algorithm A , we define $A^{(1)}$ as the algorithm that *fakes* its first query, in the following sense: If the first query made by A to the oracle is some interval I , then the algorithm $A^{(1)}$ does not call ICOND on I but instead chooses a point s uniformly at random from I and then behaves exactly as A would behave if the ICOND oracle had returned s in response to the query I . More generally, we define $A^{(k)}$ for all $k \geq 0$ as the algorithm behaving like A but faking its first k queries (note that $A^{(0)} = A$).

Let $\mathfrak{A}^{(k),N}$ denote the distribution over transcripts (of length q) when the distribution D is drawn from \mathcal{P}_{No} and the algorithm is $A^{(k)}$. Note that by this definition we have $\mathfrak{A}^{(0),N} = \mathfrak{A}^N$ and $\mathfrak{A}^{(q),N} = \mathfrak{A}^Y$.¹⁰ For a fixed outcome Δ of the block size and y of the offset in the construction of $D \sim \mathcal{P}_{\text{No}}$, we write $\mathfrak{A}_{(\Delta,y)}^{(k),N}$ to denote distribution $\mathfrak{A}^{(k),N}$ conditioned on that particular outcome of the block size and offset.

As in the non-adaptive case, in order to prove Theorem 18, it is sufficient to prove that the transcripts for uniform and “No”-distributions are close in total variation distance; i.e, that

$$d_{\text{TV}}(\mathfrak{A}^Y, \mathfrak{A}^N) \leq 1/4. \quad (81)$$

The key lemma used to prove this is the following lemma, which bounds the variation distance between the transcripts of $A^{(k)}$ (the variant of A that fakes its first k queries) and $A^{(k+1)}$, conditioned on a high-probability event over the choice of block size and offset.

Lemma 33 *For all $k \geq 0$, with probability at least $1 - \eta(N)$ over the choice of (Δ, y) in the draw of D from \mathcal{P}_{No} , we have*

$$d_{\text{TV}}\left(\mathfrak{A}_{(\Delta,y)}^{(k),N}, \mathfrak{A}_{(\Delta,y)}^{(k+1),N}\right) = \beta(k, N) \quad (82)$$

where $\eta(N) = O\left(\frac{\log \log N}{\log N}\right)$ and $\beta(k, N) = O\left(\frac{k}{(\log N)^2}\right)$.

To see why Equation((81)) follows from this lemma, first observe that

$$d_{\text{TV}}(\mathfrak{A}^N, \mathfrak{A}^Y) = d_{\text{TV}}\left(\mathfrak{A}^{(0),N}, \mathfrak{A}^{(q),N}\right) \leq \sum_{k=0}^{q-1} d_{\text{TV}}\left(\mathfrak{A}^{(k),N}, \mathfrak{A}^{(k+1),N}\right).$$

Let us say that a pair (Δ, y) for which Equation (82) holds is *good*. We now require the following variant of Fact 32:

Fact 34 *Let D_1, D_2 be two distributions over the same finite set X . Let E be an event such that $D_i[E] = \alpha_i \leq \alpha$ for $i = 1, 2$ and the conditional distributions $(D_1)_{\bar{E}}$ and $(D_2)_{\bar{E}}$ are statistically close, i.e. $d_{\text{TV}}((D_1)_{\bar{E}}, (D_2)_{\bar{E}}) = \beta$. Then $d_{\text{TV}}(D_1, D_2) \leq \alpha + \beta$.*

¹⁰We observe that there is no point in defining the corresponding distribution over transcripts for the “yes” case, as faking a call is exactly the same as querying the oracle when the underlying distribution is uniform.

Proof: As in the proof of Fact 32, let us write

$$2d_{\text{TV}}(D_1, D_2) = \sum_{x \in X \setminus E} |D_1(x) - D_2(x)| + \sum_{x \in E} |D_1(x) - D_2(x)|.$$

We may upper bound $\sum_{x \in E} |D_1(x) - D_2(x)|$ by $\sum_{x \in E} (D_1(x) + D_2(x)) = D_1(E) + D_2(E) = \alpha_1 + \alpha_2$; furthermore,

$$\begin{aligned} \sum_{x \in \bar{E}} |D_1(x) - D_2(x)| &= \sum_{x \in \bar{E}} |(D_1)_{\bar{E}}(x) \cdot D_1(\bar{E}) - (D_2)_{\bar{E}}(x) \cdot D_2(\bar{E})| \\ &\leq D_1(\bar{E}) \cdot \sum_{x \in \bar{E}} |(D_1)_{\bar{E}}(x) - (D_2)_{\bar{E}}(x)| + |D_1(\bar{E}) - D_2(\bar{E})| \cdot (D_2)_{\bar{E}}(\bar{E}) \\ &\leq (1 - \alpha_1) \cdot (2\beta) + |\alpha_2 - \alpha_1| \cdot 1 \leq 2\beta + |\alpha_2 - \alpha_1| \end{aligned}$$

Thus $2d_{\text{TV}}(D_1, D_2) \leq 2\beta + |\alpha_1 - \alpha_2| + \alpha_1 + \alpha_2 = 2\beta + 2 \max\{\alpha_1, \alpha_2\} \leq 2(\alpha + \beta)$, and the fact is established. ■

Applying this fact to distributions $\mathfrak{A}^{(k),N}$ and $\mathfrak{A}^{(k+1),N}$ and the event $E = \{(\Delta, y) \text{ is good}\}$, we get that $d_{\text{TV}}(\mathfrak{A}^{(k),N}, \mathfrak{A}^{(k+1),N}) \leq \eta(N) + \beta(k, N)$, and consequently we have $d_{\text{TV}}(\mathfrak{A}^N, \mathfrak{A}^Y) \leq \sum_{k=0}^q (\eta(N) + \beta(k, N))$, which is at most $1/4$ for a suitable choice of the absolute constant τ .

Thus it remains only to prove Lemma 33. In the next subsection we describe an alternative view of the random draw of $D \sim \mathcal{P}_{\text{No}}$, and in the following subsection we use this alternate view to prove the lemma.

8.2.4 Extended transcripts and drawing $D \sim \mathcal{P}_{\text{No}}$ on the fly.

Observe that the testing algorithm, seeing only pairs of queries and answers, does not have direct access to all the underlying information – namely, in the case of a “No”-distribution, whether the part of the block the sample comes from is high or low. It will be useful for us to consider an “extended” version of the transcripts, which includes this information though it is not directly available to the algorithm.

Definition 8 *With the same notation as in Definition 7, the extended transcript of a sequence of queries made by A and the corresponding responses is a sequence $\mathcal{E} = (I_i, s_i, b_i)_{i \in [q]}$ of triples, where I_i and s_i are as before, and $b_i \in \{\downarrow, \uparrow, \uparrow\downarrow\}$ is the profile of the block s_i belongs to. We define $\mathcal{E}|_k$ to be the length- k prefix of an extended transcript \mathcal{E} .*

We also consider extended transcripts of algorithms that fake queries. For any $0 \leq k \leq q$ we let $\mathfrak{E}^{(k),N}$ denote the distribution over extended transcripts when $D \sim \mathcal{P}_{\text{No}}$ and the algorithm is $A^{(k)}$. We may think of a draw from $\mathfrak{E}^{(k),N}$ in the following way: every time a query I_i from the algorithm is answered with a point s_i (either answered from the oracle or faked), the corresponding block profile b_i of s_i is “unveiled” and included as the third component of that triple in the extended transcript. (We go into significantly more detail on this below.) Observe that if we discard the third component b_i of each triple, then $\mathfrak{E}^{(0),N}$ corresponds precisely to the distribution \mathfrak{A}^N of regular transcripts and $\mathfrak{E}^{(k),N}$ corresponds precisely to the distribution \mathfrak{A}^Y of regular transcripts. Similar to

before, for a fixed outcome Δ of the block size and y of the offset in the construction of $D \sim \mathcal{P}_{\text{No}}$, we write $\mathfrak{E}_{(\Delta, y)}^{(k), N}$ to denote distribution $\mathfrak{E}^{(k), N}$ conditioned on that particular outcome of the block size and offset.

Our proof of Lemma 33 takes advantage of the fact that one can view the draw of a “No”-distribution from \mathcal{P}_{No} as being done “on the fly” during the course of algorithm A ’s (or, more generally, $A^{(k)}$ ’s) execution. More precisely, the size Δ and the offset y are drawn at the very beginning, but we may view the profile vector ϑ as having its components chosen independently, coordinate by coordinate, only as A interacts with **ICOND** – each time an element s_i is obtained in response to the i -th query, only then is the corresponding element b_i of the profile vector ϑ chosen (if it was not already determined by previous calls to **ICOND**). We now describe how the coordinates of the profile vector ϑ are generated sequentially as A interacts with **ICOND**.

Consider the ℓ -th query I_ℓ that A makes to **ICOND** $_D$. Inductively some coordinates of ϑ may have been already set by previous queries. Let B_1, \dots, B_k be the blocks that I_ℓ intersects.

1. If both of the outermost blocks B_1, B_k have had their bits in ϑ set already by previous queries, then these settings ($\downarrow\uparrow$ or $\uparrow\downarrow$) completely determine the probability under **ICOND** $_D$ that each block B_1, \dots, B_k is the block from which s_ℓ will be chosen.¹¹ The algorithm draws a block $B_i \in \{B_1, \dots, B_k\}$ according to these probabilities. If the coordinate of ϑ corresponding to B_i has already been determined by a previous query, then s_i is drawn from the correct distribution (as determined by the $\downarrow\uparrow$ or $\uparrow\downarrow$ setting) over B_i , and b_i is set accordingly. Otherwise a fair coin is tossed, b_ℓ is set either to $\downarrow\uparrow$ or to $\uparrow\downarrow$ depending on the outcome, s_ℓ is then drawn from the correct distribution (as determined by b_ℓ) and the corresponding coordinate of ϑ is set to b_i . The triple (I_ℓ, s_ℓ, b_ℓ) is taken as the ℓ -th element of the extended transcript.
2. If either of the outermost blocks B_1, B_k (or both) have not had their bits in ϑ set by previous queries, then again the probability under **ICOND** $_D$ of having s_ℓ belong to each block B_1, \dots, B_k is completely determined (the probability allocated to the “undetermined” outermost blocks is exactly proportionate to the amount of the block contained in I_ℓ – since the block is still undetermined, the probabilities average out to uniform). The algorithm then proceeds as above, using the correct probabilities in this case – it draws a block $B_i \in \{B_1, \dots, B_k\}$ according to these probabilities and either sets a new coordinate of ϑ if necessary or uses the old coordinate, draws s_ℓ from the block B_i , and returns s_ℓ, b_ℓ as above. As before the triple (I_ℓ, s_ℓ, b_ℓ) is taken as the ℓ -th element of the extended transcript.

This completes our description of how one may view the coordinates of the profile vector ϑ as being generated sequentially as A interacts with **ICOND**.

We next observe that we can also adopt such a view for an algorithm $A^{(k)}$, which fakes the first k queries. Recall that when an algorithm $A^{(k)}$ fakes one of the first k queries (say the ℓ -th query) it draws s_ℓ uniformly at random in I_ℓ . In the corresponding extended transcript $\mathcal{E}|_k$ of $A^{(k)}$ the profile bits are set even though $A^{(k)}$ has no access to them; we now describe how these bits

¹¹We could give an exact expression for these probabilities but it would be cumbersome and is not necessary for the rest of the argument – all we need is that the probabilities are well defined. Note that any two “internal” blocks B_2, \dots, B_{k-1} have equal probability regardless of whether or not their coordinates of ϑ have yet been set.

are set. Consider the ℓ -th query I_ℓ ($\ell \leq k$) that $A^{(k)}$ fakes. Let B_1, \dots, B_j be the blocks that I_ℓ intersects. The point s_ℓ is chosen uniformly at random from I_ℓ , the bit b_ℓ is set as follows, and the triple (I_ℓ, s_ℓ, b_ℓ) is taken as the ℓ -th element of the extended transcript:

1. If s_ℓ belongs to a block whose profile was already set during the interaction for the ℓ' -th query for some $\ell' < \ell$, then the bit b_ℓ is set to $b_\ell = b_{\ell'}$.
2. If I_ℓ is completely included within a single half-block, then b_ℓ is set uniformly at random to one of $\{\downarrow\uparrow, \uparrow\downarrow\}$.
3. If s_ℓ belongs to a block B that is completely included in I_ℓ , then we look at the half of the block s_ℓ belongs to, and toss a biased coin to set its profile $b_\ell \in \{\downarrow\uparrow, \uparrow\downarrow\}$: If s_ℓ belongs to the first half, then the coin toss's probabilities are $((1 - 2\epsilon)/2, (1 + 2\epsilon)/2)$; otherwise, they are $((1 + 2\epsilon)/2, (1 - 2\epsilon)/2)$.
4. if s_ℓ belongs to one of the two end blocks, then:
 - (a) if I_ℓ contains at most half of the block that s_ℓ belongs to, again b_ℓ is set to $\downarrow\uparrow$ or $\uparrow\downarrow$ with equal probability;
 - (b) if I_ℓ contains a fraction $\frac{1}{2} + x$ of the block that s_ℓ belongs to, then the coin is biased, with probability for setting the part of the block in which s_ℓ lies to “high” either $\frac{1}{1 + \frac{\Delta}{2x} \cdot \frac{1-2\epsilon}{1+2\epsilon}}$ if this part is the small one (the x portion), and $\frac{1}{1 + \frac{2x}{\Delta} \cdot \frac{1-2\epsilon}{1+2\epsilon}}$ if s_ℓ is in the complete half-block.

It follows from the foregoing description that in both cases (an algorithm A that does not fake queries, or an algorithm $A^{(k)}$ that does), the process described above for choosing the profile vector of D “on the fly” indeed corresponds to drawing D from \mathcal{P}_{No} . Similarly, the resulting distributions over extended transcripts correspond to draws from $\mathfrak{E}^{(0),N}$ (in the A case) and from $\mathfrak{E}^{(k),N}$ (in the $A^{(k)}$ case) respectively.

8.2.5 Proof of Lemma 33

Fix $0 \leq k \leq q$ and recall the definitions of $A^{(k)}$, $A^{(k+1)}$ and $\mathfrak{E}_{(\Delta,y)}^{(k),N}$, $\mathfrak{E}_{(\Delta,y)}^{(k+1),N}$. (For the sake of concision, we shall write $\mathcal{E}^{(k)}$ and $\mathcal{E}^{(k+1)}$ for the transcript random variables, and $\mathfrak{E}^{(k)}$, $\mathfrak{E}^{(k+1)}$ for their distribution.)

We first observe that prior to the $(k+1)$ st query, the distribution of length- k transcript prefixes is exactly the same under $\mathfrak{E}^{(k)}$ and $\mathfrak{E}^{(k+1)}$, so we have $d_{\text{TV}}(\mathfrak{E}^{(k)}|_k, \mathfrak{E}^{(k+1)}|_k) = 0$. For any fixed extended transcript $E = (E_1, \dots, E_q)$, we thus have $p_k(E) \stackrel{\text{def}}{=} \Pr[\mathcal{E}^{(k)}|_k = E|_k] = \Pr[\mathcal{E}^{(k+1)}|_k = E|_k]$, and

$$\begin{aligned} & \left| \Pr[\mathcal{E}^{(k)} = E] - \Pr[\mathcal{E}^{(k+1)} = E] \right| = \\ & \underbrace{\left| \Pr[\mathcal{E}_{k+1}^{(k)} = E_{k+1} \mid \mathcal{E}^{(k)}|_k = E|_k] - \Pr[\mathcal{E}_{k+1}^{(k+1)} = E_{k+1} \mid \mathcal{E}^{(k+1)}|_k = E|_k] \right|}_{(i)} \cdot p_k(E). \end{aligned}$$

Consequently it suffices to only analyze the distribution of the $(k + 1)$ -st component $(I_{k+1}, s_{k+1}, b_{k+1})$ of the transcript random variables $\mathcal{E}^{(k)}, \mathcal{E}^{(k+1)}$. Because both variants $A^{(k)}$ and $A^{(k+1)}$ of the algorithm select the query interval I_{k+1} based only on the length- k transcript prefix (identical in both cases, since we conditioned on this in (\dagger)) and their internal randomness (identically distributed in both cases), both algorithms will query any given interval I_{k+1} with the same probability, so we only have to deal with the remainder of the triple, namely the last two entries (s_{k+1}, b_{k+1}) .

Fix any given outcome I'_{k+1} for I_{k+1} . We now define a set G of possible outcomes (Δ, y) for the block size and the offset in the draw of a “No”-distribution $D \sim \mathcal{P}_{\text{No}}$ as follows. The outcomes (Δ, y) that are *not* in G are those that fall into one of the following three categories (the first two corresponding to two of the “bad events” from our earlier discussion):

- (a) $\Delta / \log N \leq |I'_{k+1}| \leq \Delta \cdot (\log N)^2$;
- (b) $|I'_{k+1}| < \Delta / \log N$ and I'_{k+1} intersects two blocks;
- (c) $|I'_{k+1}| < \Delta / \log N$ and I'_{k+1} is contained entirely within a single block but not entirely within a half block (it intersects both the first half and second half of the block it is in).

We claim that $\Pr[(\Delta, y) \in G] \geq 1 - \eta(N)$ (where $\eta(N) = O(\frac{\log \log N}{\log N})$ is as defined in Lemma 33). This is because the probability of (a) (over a random choice of Δ) is at most $O(\frac{\log \log N}{\log N})$ and the probability of both (b) and (c) (over a random choice of y , for any fixed outcome of Δ satisfying $|I'_{k+1}| < \Delta / \log N$) is at most $O(1 / \log N)$, by applying an argument that is essentially the same as the one used to argue Items (i) and (ii) in Section 8.2.2.

Fix any pair $(\Delta, y) \in G$. It remains only to argue that Equation (82) holds for such a pair; from the discussion above, it is enough to consider outcomes I'_{k+1} for I_{k+1} that fall into one of the following categories:

- (d) $|I'_{k+1}| < \Delta / \log N$ and I'_{k+1} is contained entirely within a single half block (either the first half of a block or the second half of a block);
- (e) $|I'_{k+1}| > \Delta \cdot (\log N)^2$.

Fix an outcome of $|I'_{k+1}|$ that satisfies (d). We claim that the variation distance between (the distribution of (s_{k+1}, b_{k+1}) under $\mathfrak{E}^{(k)}$) and (the distribution of (s_{k+1}, b_{k+1}) under $\mathfrak{E}^{(k+1)}$) is zero. This is because as described in the previous subsection, in both situations ($A^{(k)}$ or $A^{(k+1)}$) the sample s is uniformly selected from I_{k+1} , and b_{k+1} is either set to the previously determined value (if a profile for this block was previously determined in an earlier step) or to a randomly selected element of $\{\downarrow, \uparrow\}$ otherwise (see items 1 and 2 p.76). Thus in Case (d) both distributions $\mathfrak{E}^{(k)}$ and $\mathfrak{E}^{(k+1)}$ over the $(k + 1)$ -th triple $(I_{k+1}, s_{k+1}, b_{k+1})$ are identical, conditioned on $I_{k+1} = I'_{k+1}$.

Now fix an outcome of $|I'_{k+1}|$ that satisfies (e), so $|I'_{k+1}| > \Delta \cdot (\log N)^2$. Intuitively, the probability that in either setting ($A^{(k)}$ or $A^{(k+1)}$) the point s_{k+1} belongs to either one of the “already touched” blocks (blocks that contain some previous sample s_ℓ , $\ell \leq k$), or to one of the outermost blocks that

I_{k+1} overlaps but does not fully contain, is very small. (In more detail, an analysis that is essentially the same as that of (iii) and (iv) in Section 8.2.2 gives that the random choice of s_{k+1} (in either setting) hits a “previously touched” block or one of the two outermost blocks with probability at most $O(k/(\log N)^2)$.) Conditioned on this not occurring, in both settings s_{k+1} is uniformly distributed among all other blocks, and the discussion of the previous subsection implies that b_{k+1} is high (respectively, low) with probability $\frac{1+2\epsilon}{2}$ ($\frac{1-2\epsilon}{2}$, respectively) depending on the half-block s_{i+1} falls into, so the variation distance between the two distributions over (s_{k+1}, b_{k+1}) under this conditioning is zero. Applying Fact 32, we get that the two distributions over triples $(I'_{k+1}, s_{k+1}, b_{k+1})$ have variation distance at most $O(k/(\log N)^2)$. Averaging over all possible outcomes of I'_{k+1} , we get that the two distributions over the $(k+1)$ -st component $(I_{k+1}, s_{k+1}, b_{k+1})$ have variation distance at most $\beta(k, N) = O(k/(\log N)^2)$. This establishes Lemma 33 and concludes the proof of the theorem. ■

9 Conclusion

We have introduced a new conditional sampling framework for testing probability distributions and shown that it allows significantly more query-efficient algorithms than the standard framework for a range of problems. This new framework presents many potential directions for future work.

One specific goal is to strengthen the upper and lower bounds for problems studied in this paper. As a concrete question along these lines, we conjecture that COND algorithms for testing equality of two unknown distributions D_1 and D_2 over $[N]$ require $(\log N)^{\Omega(1)}$ queries. A broader goal is to study more properties of distributions beyond those considered in this paper; natural candidates here, which have been well-studied in the standard model, are monotonicity (for which we have preliminary results), independence between marginals of a joint distribution, and entropy. Yet another goal is to study distributions over other structured domains such as the Boolean hypercube $\{0, 1\}^n$ – here it would seem natural to consider “subcube” queries, analogous to the ICOND queries we considered when the structured domain is the linearly ordered set $[N]$. A final broad goal is to study distribution *learning* (rather than testing) problems in the conditional sampling framework.

References

- [AJ06] J. A. Adell and P. Jodra. Exact kolmogorov and total variation distances between some familiar discrete distributions. *J. of Ineq. App.*, 2006(1):64307, 2006.
- [BDKR05] T. Batu, S. Dasgupta, R. Kumar, and R. Rubinfeld. The complexity of approximating the entropy. *SICOMP*, 35(1):132–150, 2005.
- [BFF⁺01] T. Batu, E. Fischer, L. Fortnow, R. Kumar, R. Rubinfeld, and P. White. Testing random variables for independence and identity. In *Proceedings of FOCS*, pages 442–451, 2001.
- [BFR⁺00] T. Batu, L. Fortnow, R. Rubinfeld, W. D. Smith, and P. White. Testing that distributions are close. In *Proceedings of FOCS*, pages 189–197, 2000.

- [BFR⁺10] T. Batu, L. Fortnow, R. Rubinfeld, W. D. Smith, and P. White. Testing closeness of discrete distributions. Technical Report abs/1009.5397, 2010. This is a long version of [BFR⁺00].
- [BFRV11] A. Bhattacharyya, E. Fischer, R. Rubinfeld, and P. Valiant. Testing monotonicity of distributions over general partial orders. In *Proceedings of ITCS*, pages 239–252, 2011.
- [BKR04] T. Batu, R. Kumar, and R. Rubinfeld. Sublinear algorithms for testing monotone and unimodal distributions. In *Proceedings of STOC*, pages 381–390, 2004.
- [CFGM13] S. Chakraborty, E. Fischer, Y. Goldhirsh, and A. Matsliah. On the power of conditional samples in distribution testing. In *Proceedings of ITCS*, 2013. To appear.
- [DDS⁺13] C. Daskalakis, I. Diakonikolas, R. Servedio, G. Valiant, and P. Valiant. Testing k -modal distributions: Optimal algorithms via reductions. In *Proceedings of SODA*, 2013.
- [DP09] D. Dubhashi and A. Panconesi. *Concentration of measure for the analysis of randomized algorithms*. Cambridge University Press, Cambridge, 2009.
- [Fis01] E. Fischer. The art of uninformed decisions: A primer to property testing. *BEATCS*, 75:97–126, 2001.
- [GGR98] O. Goldreich, S. Goldwasser, and D. Ron. Property testing and its connection to learning and approximation. *JACM*, 45(4):653–750, 1998.
- [Gol10] O. Goldreich, editor. *Property Testing: Current Research and Surveys*. Springer, 2010. LNCS 6390.
- [GR00] O. Goldreich and D. Ron. On testing expansion in bounded-degree graphs. Technical Report TR00-020, ECCC, 2000.
- [ILR12] P. Indyk, R. Levi, and R. Rubinfeld. Approximating and Testing k -Histogram Distributions in Sub-linear Time. In *Proceedings of PODS*, pages 15–22, 2012.
- [Ma81] S. K. Ma. Calculation of entropy from data of motion. *J. Stat. Phys.*, 26(2), 1981.
- [MR95] R. Motwani and P. Raghavan. *Randomized Algorithms*. Cambridge University Press, New York, NY, 1995.
- [Pan04] L. Paninski. Estimating entropy on m bins given fewer than m samples. *IEEE-IT*, 50(9):2200–2203, 2004.
- [Pan08] L. Paninski. A coincidence-based test for uniformity given very sparsely sampled discrete data. *IEEE-IT*, 54(10):4750–4755, 2008.
- [Rey11] L. Reyzin. Extractors and the leftover hash lemma. <http://www.cs.bu.edu/reyzin/teaching/s11cs937/notes-leo-1.pdf>, March 2011.
- [Ron08] D. Ron. Property Testing: A Learning Theory Perspective. *FnTML*, 1(3):307–402, 2008.
- [Ron10] D. Ron. Algorithmic and analysis techniques in property testing. *FnTCS*, 5:73–205, 2010.

- [RRSS09] S. Raskhodnikova, D. Ron, A. Shpilka, and A. Smith. Strong lower bounds for approximating distributions support size and the distinct elements problem. *SICOMP*, 39(3):813–842, 2009.
- [RS96] R. Rubinfeld and M. Sudan. Robust characterization of polynomials with applications to program testing. *SICOMP*, 25(2):252–271, 1996.
- [RS09] R. Rubinfeld and R. A. Servedio. Testing monotone high-dimensional distributions. *RSA*, 34(1):24–44, January 2009.
- [Val11] P. Valiant. Testing symmetric properties of distributions. *SICOMP*, 40(6):1927–1968, 2011.
- [VV10a] G. Valiant and P. Valiant. A CLT and tight lower bounds for estimating entropy. Technical Report TR10-179, ECCC, 2010.
- [VV10b] G. Valiant and P. Valiant. Estimating the unseen: A sublinear-sample canonical estimator of distributions. Technical Report TR10-180, ECCC, 2010.
- [VV11] G. Valiant and P. Valiant. Estimating the unseen: an $n/\log(n)$ -sample estimator for entropy and support size, shown optimal via new CLTs. In *Proceedings of STOC*, pages 685–694, 2011. See also [VV10a] and [VV10b].