



A Chasm Between Identity and Equivalence Testing with Conditional Queries

Jayadev Acharya* Clément L. Canonne† Gautam Kamath‡

April 18, 2015

Abstract

A recent model for property testing of probability distributions [CFGM13, CRS15] enables tremendous savings in the sample complexity of testing algorithms, by allowing them to condition the sampling on subsets of the domain.

In particular, Canonne, Ron, and Servedio [CRS15] showed that, in this setting, testing identity of an unknown distribution D (i.e., whether $D = D^*$ for an explicitly known D^*) can be done with a *constant* number of samples, independent of the support size n – in contrast to the required \sqrt{n} in the standard sampling model. However, it was unclear whether the same held for the case of testing equivalence, where *both* distributions are unknown. Indeed, while Canonne, Ron, and Servedio [CRS15] established a $\text{polylog}(n)$ -query upper bound for equivalence testing, very recently brought down to $\tilde{O}(\log \log n)$ by Falahatgar et al. [FJO⁺15], whether a dependence on the domain size n is necessary was still open, and explicitly posed by Fischer at the Bertinoro Workshop on Sublinear Algorithms [Sublinear.info, Problem 66]. In this work, we answer the question in the positive, showing that any testing algorithm for equivalence must make $\Omega(\sqrt{\log \log n})$ queries in the conditional sampling model. Interestingly, this demonstrates an intrinsic qualitative gap between identity and equivalence testing, absent in the standard sampling model (where both problems have sampling complexity $n^{\Theta(1)}$).

Turning to another question, we investigate the complexity of support size estimation. We provide a doubly-logarithmic upper bound for the adaptive version of this problem, generalizing work of Ron and Tsur [RT14] to our weaker model. We also establish a logarithmic lower bound for the non-adaptive version of this problem. This latter result carries on to the related problem of non-adaptive uniformity testing, an exponential improvement over previous results that resolves an open question of Chakraborty, Fischer, Goldhirsh, and Matsliah [CFGM13].

*EECS, MIT. Email: jayadev@csail.mit.edu. Research supported by grant from MITEI-Shell program.

†Columbia University. Email: cannonne@cs.columbia.edu. Research supported by NSF CCF-1115703 and NSF CCF-1319788.

‡EECS, MIT. Email: g@csail.mit.edu.

1 Introduction

“No, Virginia, there is no constant-query tester.”

Understanding properties and characteristics of an unknown probability distribution is a fundamental problem in statistics, and one that has been thoroughly studied. However, it is only since the seminal work of Goldreich and Ron [GR00] and Batu et al. [BFR⁺00] that the problem has been considered through the lens of theoretical computer science, more particularly in the setting of *property testing*.

Over the following decade, a flurry of subsequent work explored and delved into this new area, resulting in a better and often complete understanding of a number of questions in distributional property testing (see e.g. [GR00, BFF⁺01, BKR04, Pan08, RS09, ADJ⁺11, BFRV11, Rub12, ILR12, ADJ⁺12, CDVV14, VV14] or [Can15] for a survey). In many cases, these culminated in provably sample-optimal algorithms. However, the standard setting of distribution testing, where one only obtains independent samples from an unknown distribution D , does not encompass *all* scenarios one may encounter. In recent years, stronger models have thus been proposed to capture more specific situations [GMV06, CFGM13, CRS15, LRR13, CR14]: among these is the *conditional oracle model* [CFGM13, CRS15] which will be the focus of our work. In this setting, the testing algorithms are given the ability to sample from conditional distributions: that is, to specify a subset S of the domain and obtain samples from D_S , the distribution induced by D on S (the formal definition of the model can be found in [Definition 2.1](#)). In particular, the hope is that allowing algorithms to have stronger interactions with the unknown underlying distributions might significantly reduce the number of samples they need, thereby sidestepping the strong lower bounds that hold in the standard sampling model.

1.1 Background and previous work

We focus in this paper on proving lower bounds for testing two extremely natural properties of distributions, namely *equivalence testing* (“are these two datasets identically distributed?”) and *support size estimation* (“how many different outcomes can actually be observed?”). Along the way, we use some of the techniques we develop to obtain an upper bound on the query complexity of the latter. We state below the informal definition of these two problems, along with closely related ones (uniformity and identity testing). Hereafter, “oracle access” to a distribution D over $[n] = \{1, \dots, n\}$ means access to samples generated independently from D .

Uniformity testing: granted oracle access to D , decide whether $D = \mathcal{U}$ (the uniform distribution on $[n]$) or is far from it;

Identity testing: granted oracle access to D and the full description of a fixed D^* , decide whether $D = D^*$ or is far from it;

Equivalence (closeness) testing: granted independent oracle accesses to D_1, D_2 (both unknown), decide whether $D_1 = D_2$ or D_1, D_2 are far from each other.

Support size estimation: granted oracle access to D , output an estimate of the size of the support¹ $\text{supp}(D) = \{x : D(x) > 0\}$, accurate within a multiplicative factor.

¹For this problem, it is typically assumed that all points in the support have probability mass at least $\Omega(1)/n$, as without such guarantee it becomes impossible to give any non-trivial estimate (consider for instance a distribution D such that $D(i) \propto 1/2^{in}$).

It is not difficult to see that each of the first three problems generalizes the previous, and is therefore at least as hard. All of these tasks are known to require sample complexity $n^{\Omega(1)}$ in the standard sampling model (SAMP); yet, as prior work [CFG13, CRS15] shows, their complexity decreases tremendously when one allows the stronger type of access to the distribution(s) provided by a conditional sampling oracle (COND). For the problems of uniformity testing and identity testing, the sample complexity even becomes a constant provided the testing algorithm is allowed to be *adaptive* (i.e. when the next queries it makes can depend on the samples it previously obtained).

Testing uniformity and identity. Given the complete description of a distribution D^* over $[n]$, a parameter $\varepsilon > 0$, and oracle access to a distribution D , identity testing asks to distinguish the case $D_1 = D^*$ from where their total variation distance $d_{\text{TV}}(D, D^*)$ is at least ε . This is a generalization of uniformity testing, where D^* is taken to be the uniform distribution over $[n]$. The complexity of these tasks is well-understood in the sampling model; in particular, it is known that for both uniformity and identity testing $\Theta(\sqrt{n}/\varepsilon^2)$ samples are necessary and sufficient (see [GR00, BFR⁺10, Pan08, VV14] for the tight bounds on these problems).

The uniformity testing problem emphasizes the additional flexibility granted by conditional sampling: as Canonne, Ron, and Servedio [CRS15] showed, in this setting only $\tilde{O}(1/\varepsilon^2)$ adaptive queries now suffice (and this is optimal, up to logarithmic factors). They further prove that identity testing has constant sample complexity as well, namely $\tilde{O}(1/\varepsilon^4)$ – very recently improved to a near-optimal $\tilde{O}(1/\varepsilon^2)$ by Falahatgar et al. [FJO⁺15]. The power of the COND model is evident from the fact that a task requiring polynomially many samples in the standard model can now be achieved with a number of samples *independent of the domain size n* .

Focusing on the case of non-adaptive algorithms, Chakraborty et al. [CFG13] describe a $\text{poly}(\log n, 1/\varepsilon)$ -query tester for uniformity, showing that even without the full power of conditional queries one can still get an exponential improvement over the standard sampling setting. They also obtain an $\Omega(\log \log n)$ lower bound for this problem, and leave open the possibility of improving this lower bound up to a logarithmic dependence. The present work answers this question, establishing that any non-adaptive uniformity tester must perform $\Omega(\log n)$ conditional queries.

Testing equivalence. A natural generalization of these two testing problems is the question of *equivalence testing*, defined as follows. Given oracle access to two unknown distributions D_1 and D_2 over $[n]$ and a parameter $\varepsilon > 0$, equivalence testing asks to distinguish between the cases $D_1 = D_2$ and $d_{\text{TV}}(D_1, D_2) > \varepsilon$. This problem has been extensively studied over the past decade, and its sample complexity is now known to be $\Theta(\max(n^{2/3}/\varepsilon^{4/3}, \sqrt{n}/\varepsilon^2))$ in the sampling model [BFR⁺10, Val11, CDVV14].

In the COND setting, Canonne, Ron, and Servedio showed that equivalence testing is possible with only $\text{poly}(\log n, 1/\varepsilon)$ queries. Concurrent to our work, Falahatgar et al. [FJO⁺15] brought this upper bound down to $\tilde{O}((\log \log n)/\varepsilon^5)$, a *doubly exponential* improvement over the $n^{\Omega(1)}$ samples needed in the standard sampling model. However, these results still left open the possibility of a constant query complexity: given that both uniformity and identity testing admit constant-query testers, it is natural to wonder where equivalence testing lies².

²It is worth noting that an $\Omega(\log^c n)$ lower bound was known for equivalence testing in a weaker version of the conditional oracle, PAIRCOND (where the tester’s queries are restricted to being either $[n]$ or subsets of size 2 [CRS15]).

This question was explicitly posed by Fischer at the Bertinoro Workshop on Sublinear Algorithms 2014 [Sublinear.info, Problem 66]: in this paper, we make decisive progress in answering it, ruling out the possibility of any constant-query tester for equivalence. Along with the upper bound of Falahatgar et al. [FJO⁺15], our results essentially settle the dependence on the domain size, showing that $(\log \log n)^{\Theta(1)}$ samples are both necessary and sufficient.

Support size estimation. Finally, the question of approximating the support size of a distribution has been considered by Raskhodnikova et al. [RRSS09], where it was shown that obtaining *additive* estimates requires sample complexity almost linear in n . Subsequent work by Valiant and Valiant [VV11, VV10a] settles the question, establishing that an $n/\log n$ dependence is both necessary and sufficient. Note that the proof of their lower bound translates to multiplicative approximations as well, as they rely on the hardness of distinguishing a distribution with support $s \leq n$ from a distribution with support $s + \varepsilon n \geq (1 + \varepsilon)s$. To the best of our knowledge, the question of getting a multiplicative-factor estimate of the support size of a distribution given conditional sampling access has not been previously considered. We provide upper and lower bounds for both the adaptive and non-adaptive versions of this problem.

1.2 Our results

In this work, we make significant progress in each of the problems introduced in the previous section, yielding a better understanding of their intrinsic query complexities. We prove four results pertaining to the sample complexity of equivalence testing, support size estimation, and uniformity testing in the COND framework.

Problem	COND model	Standard model
Are D_1, D_2 (both unknown) equivalent? (adaptive)	$\tilde{O}\left(\frac{\log \log n}{\varepsilon^5}\right)$ [FJO ⁺ 15] $\Omega(\sqrt{\log \log n})$ [this work]	$\Theta\left(\max\left(\frac{n^{2/3}}{\varepsilon^{4/3}}, \frac{n^{1/2}}{\varepsilon^2}\right)\right)$ [CDVV14]
What is the support size of D ? (adaptive)	$\tilde{O}\left(\frac{\log \log n}{\varepsilon^3}\right)$ [this work] $\Omega(\sqrt{\log \log n})$ [CFGM13] (†)	$\Theta\left(\frac{n}{\log n}\right)$ [VV10a]
What is the support size of D ? (non-adaptive)	$O(\text{poly}(\log n, 1/\varepsilon))$ [this work] $\Omega(\log n)$ [this work]	
Is D uniform over the domain? (non-adaptive)	$\tilde{O}\left(\frac{\log^5 n}{\varepsilon^6}\right)$ [CFGM13] $\Omega(\log n)$ [this work]	$\Theta\left(\frac{\sqrt{n}}{\varepsilon^2}\right)$ [GR00, BFR ⁺ 10, Pan08]

Table 1: Summary of results. Note that the lower bound (†) can also be easily derived from our lower bound on testing equivalence.

Our main result considers the sample complexity of testing equivalence with adaptive queries under the COND model, resolving in the negative the question of whether constant-query complexity was achievable [Sublinear.info, Problem 66]. More precisely, we prove the following theorem:

Theorem 1.1 (Testing Equivalence). *Any adaptive algorithm which, given COND access to unknown distributions D_1, D_2 on $[n]$, distinguishes with probability at least $2/3$ between (a) $D_1 = D_2$ and (b) $d_{TV}(D_1, D_2) \geq \frac{1}{4}$, must have query complexity $\Omega(\sqrt{\log \log n})$.*

Combined with the recent $\tilde{O}(\log \log n)$ upper bound of Falahatgar et al. [FJO⁺15], this almost settles the sample complexity of this question. Furthermore, as the related task of identity testing *can* be performed with a constant number of queries in the conditional sampling model, this demonstrates an intriguing and intrinsic difference between the two problems. Our result can also be interpreted as showing a fundamental distinction from the usual sampling model, where both identity and equivalence testing have polynomial sample complexity.

Next, we establish a logarithmic lower bound on *non-adaptive* support size estimation, for any factor larger than a fixed constant. This improves on the result of Chakraborty et al. [CFG13], which gave a doubly logarithmic lower bound for constant factor support-size estimation.

Theorem 1.2 (Non-Adaptive Support Size Estimation). *Any non-adaptive algorithm which, given COND access to an unknown distribution D on $[n]$, estimates the size of its support up to a factor γ must have query complexity $\Omega\left(\frac{\log n}{\log \gamma}\right)$, for any $\gamma \geq \sqrt{2}$.*

Moreover, the approach used to prove this theorem also implies an analogous lower bound on *non-adaptive* uniformity testing in the conditional model, answering a conjecture of Chakraborty et al. [CFG13]:

Theorem 1.3 (Non-Adaptive Uniformity Testing). *Any non-adaptive algorithm which, given COND access to an unknown distribution D on $[n]$, distinguishes with probability at least $2/3$ between (a) $D = \mathcal{U}$ and (b) $d_{\text{TV}}(D, \mathcal{U}) \geq \frac{1}{4}$, must have query complexity $\Omega(\log n)$.*

We note that these results complement $\text{polylog}(n)$ -query upper bounds, the former of which we sketch in this paper, and the latter obtained by Chakraborty et al. [CFG13]. This shows that both of these problems have query complexity $\log^{\Theta(1)} n$ in the non-adaptive case.

Finally, we conclude with an upper bound for *adaptive* support size estimation. Specifically, we provide a $\tilde{O}(\log \log n)$ -query algorithm for support size estimation. This shows that the question becomes *double exponentially* easier when conditional samples are allowed.

Theorem 1.4 (Adaptive Support Size Estimation). *Let $\tau > 0$ be any constant. There exists an adaptive algorithm which, given COND access to an unknown distribution D on $[n]$ which has minimum non-zero probability τ/n and accuracy parameter ε makes $\tilde{O}((\log \log n)/\varepsilon^3)$ queries to the oracle and outputs a value $\tilde{\omega}$ such that the following holds. With probability at least $2/3$, $\tilde{\omega} \in [\frac{1}{1+\varepsilon} \cdot \omega, (1 + \varepsilon) \cdot \omega]$, where $\omega = |\text{supp}(D)|$.*

1.2.1 Relation to the Ron-Tsur model

Recent work of Ron and Tsur [RT14] studies a model which is slightly stronger than ours. In their setting, the algorithm still performs queries consisting of a subset of the domain. However, the algorithm is also given the promise that the distribution is uniform on a subset of the domain, and whenever a query set contains 0 probability mass the oracle explicitly indicates this is the case. Their paper provides a number of results for support size estimation in this model.

We point out two connections between our work and theirs. First, our $\Omega(\log n)$ lower bound for non-adaptive support size estimation (Theorem 1.2) leads to the same lower bound for the problem in the model of Ron and Tsur. Although lower bounds in the conditional sampling setting do not apply directly to theirs, we note that our construction and analysis still carry over. This provides a nearly tight answer to this question, which was left unanswered in their paper. Also,

our $\tilde{O}(\log \log n)$ -query algorithm for adaptive support size estimation ([Theorem 1.4](#)) can be seen as generalizing their result to the weaker conditional sampling model (most significantly, when we are not given the promise that the distribution be uniform).

1.3 Techniques and proof ideas

We now provide an overview of the techniques and arguments used to prove our results.

Lower bound on adaptive equivalence testing. In order to prove our main $\omega(1)$ lower bound on the query complexity of testing equivalence in the conditional sampling model, we have to deal with one main conceptual issue: *adaptivity*. While the standard sampling model does not, by definition, allow any choice on what the next query to the oracle should be, this is no longer the case for COND algorithms. Quantifying the power that this grants an algorithm makes things much more difficult. To handle this point, we follow the approach of Chakraborty et al. [[CFG13](#)] and focus on a restricted class of algorithms they introduce, called “core adaptive testers” (see [Section 2.2](#) for a formal definition). They show that this class of testers is equivalent to general algorithms for the purpose of testing a broad class of properties, namely those which are invariant to any permutation of the domain. Using this characterization, it remains for us to show that none of these structurally much simpler core testers can distinguish whether they are given conditional access to (a) a pair of random identical distributions (D_1, D_1) , or (b) two distributions (D_1, D_2) drawn according to a similar process, which are far apart.

At a high level, our lower bound works by designing instances where the property can be tested if and only if the support size is known to the algorithm. Our construction randomizes the support size by embedding the instance into a polynomially larger domain. Since the algorithm is only allowed a small number of queries, Yao’s Principle allows us to argue that, with high probability, a deterministic algorithm is unable to “guess” the support size. This separates queries into several cases. First, in a sense we make precise, it is somehow “predictable” whether or not a query will return an element we have previously observed. If we do, it is similarly predictable *which* element the query will return. On the other hand, if we observe a fresh element, the query set is either “too small” or “too large.” In the former case, the query will entirely miss the support, and the sampling process is identical for both types of instance. In the latter case, the query will hit a large portion of the support, and the amount of information gleaned from a single sample is minimal.

At a lower level, this process itself is reminiscent of the lower bound construction of Canonne, Ron, and Servedio [[CRS15](#)] on testing identity (with a PAIRCOND oracle), with one pivotal twist. As in their work, both D_1 and D_2 are uniform within each of $\omega(1)$ “buckets” whose size grows exponentially and are grouped into “bucket-pairs.” Then, D_2 is obtained from D_1 by internally redistributing the probability mass of each pair of buckets, so that the total mass of each pair is preserved but each particular bucket has mass going up or down by a constant factor (see [Section 3.1](#) for details of the construction). However, we now add a final step, where in both D_1 and D_2 the resulting distribution’s support is *scaled by a random factor*, effectively reducing it to a (randomly) negligible fraction of the domain. Intuitively, this last modification has the role of “blinding” the testing algorithm: we argue that unless its queries are on sets whose size somehow match (in a sense formalized in [Section 3.2](#)) this random size of the support, the sequences of samples it will obtain under D_1 and D_2 are almost identically distributed. The above discussion crucially hides many significant aspects and technical difficulties which we address in [Section 3](#). Moreover, we observe that the lower bound we obtain seems to be optimal with regard to our proofs techniques

(specifically, to the decision tree approach), and not an artifact of our lower bound instances. Namely, there appear to be conceptual barriers to strengthening our result, which would require new ideas.

Lower bound on non-adaptive support size estimation. Turning to the (non-adaptive) lower bound of [Theorem 1.2](#), we define two families of distributions \mathcal{D}_1 and \mathcal{D}_2 , where an instance is either a draw (D_1, D_2) from $\mathcal{D}_1 \times \mathcal{D}_2$, or simply (D_1, D_1) . Any distribution in \mathcal{D}_2 has support size γ times that of its corresponding distribution in \mathcal{D}_1 . Yet, we argue that no non-adaptive *deterministic* tester making too few queries can distinguish between these two cases, as the tuple of samples it will obtain from D_1 or (the corresponding) D_2 is almost identically distributed (where the randomness is over the choice of the instance itself). To show this last point, we analyze separately the case of “small” queries (conditioning on sets which turn out to be much smaller than the actual support size, and thus with high probability will not even intersect it) and the “big” ones (where the query set A is so big in front of the support size S that a uniform sample from $A \cap S$ is essentially indistinguishable from a uniform sample from A). We conclude the proof by invoking Yao’s Principle, carrying the lower bound back to the setting of non-adaptive *randomized* testers.

Interestingly, this argument essentially gives us [Theorem 1.3](#) “for free:” indeed, the big-query-set case above is handled by proving that the distribution of samples returned on those queries is indistinguishable, both for \mathcal{D}_1 and \mathcal{D}_2 , from samples obtained from the *actual* uniform distribution. Considering again the small-query-set case separately, this allows us to argue that a random distribution from (say) \mathcal{D}_1 is indistinguishable from uniform.

Upper bound on support size estimation. Our algorithm for estimating the support size to a constant factor ([Theorem 1.4](#)) is simple in spirit, and follows a guess-and-check strategy. In more detail, it first obtains a “reference point” *outside* the support, to check whether subsequent samples it may consider belong to the support. Then, it attempts to find a *rough upper bound* on the size of the support, of the form 2^{2^j} (so that only $\log \log n$ many options have to be considered); by using its reference point to check if a uniform random subset of this size contains, as it should, at least one point from the support. Once such an upper bound has been obtained using this double-exponential strategy, a refined bound is then obtained via a binary search on the new range of values for the exponent, $\{2^{j-1}, \dots, 2^j\}$. Not surprisingly, our algorithm draws on similar ideas as in [[RT14](#), [Sto85](#)], with some additional machinery to supplement the differences in the models. Interestingly, as a side-effect, this upper bound shows our analysis of [Theorem 1.1](#) to be tight up to a quadratic dependence. Indeed, the lower bound construction we consider (see [Section 3.1](#)) can be easily “defeated” if an estimate of the support size is known, and therefore cannot yield better than a $\Omega(\log \log n)$ lower bound. Similarly, this also shows that the adaptive lower bound for support size estimation of Chakraborty et al. [[CFG13](#)] is also tight up to a quadratic dependence.

Organization. The rest of the paper describes details and proofs of the results mentioned in the above discussion. In [Section 2](#), we introduce the necessary definitions and some of the tools we shall use. [Section 3](#) covers our main result on adaptive equivalence testing, [Theorem 1.1](#). In [Section 4](#) we prove our lower bounds for support size estimation and uniformity testing, and [Section 5](#) details our upper bounds for support size estimation. The reader may independently read the corresponding sections at their discretion.

2 Preliminaries

2.1 Notation and sampling models

All throughout this paper, we denote by $[n]$ the set $\{1, \dots, n\}$, and by \log the logarithm in base 2. A probability distribution over a (countable) domain $[n]$ is a non-negative function $D: [n] \rightarrow [0, 1]$ such that $\sum_{x \in [n]} D(x) = 1$. We denote by $\mathcal{U}(S)$ the uniform distribution on a set S . Given a distribution D over $[n]$ and a set $S \subseteq [n]$, we write $D(S)$ for the total probability mass $\sum_{x \in S} D(x)$ assigned to S by D . Finally, for $S \subseteq [n]$ such that $D(S) > 0$, we denote by D_S the conditional distribution of D restricted to S , that is $D_S(x) = \frac{D(x)}{D(S)}$ for $x \in S$ and $D_S(x) = 0$ otherwise.

As is usual in distribution testing, in this work the distance between two distributions D_1, D_2 on $[n]$ will be the *total variation distance*:

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

which takes value in $[0, 1]$.

In this work, we focus on the setting of *conditional access* to the distribution, as introduced and studied in [CFG13, CRS15]. We reproduce below the corresponding definition of a conditional oracle, henceforth referred to as **COND**:

Definition 2.1 (Conditional access model). Fix a distribution D over $[n]$. A **COND oracle** for D , denoted COND_D , is defined as follows: the oracle takes as input a *query set* $S \subseteq [n]$, chosen by the algorithm, 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.

Note 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 . Various definitional choices could be made to deal with this. These choices do not make significant difference in most situations, as most (adaptive) algorithms can always include in their next queries a sample previously obtained; while our lower bounds can be thought of as putting exponentially small probability mass of elements outside the support. For this reason, and for convenience, we shall hereafter assume, following Chakraborty et al., that the oracle returns in this case a sample uniformly distributed in S .

Finally, recall that a *property* \mathcal{P} of distributions over $[n]$ is a set 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 then defined as $\inf_{D' \in \mathcal{P}} d_{\text{TV}}(D, D')$. We use the standard definition of testing algorithms for properties of distributions over $[n]$, tailored for the setting of conditional access to an unknown distribution:

Definition 2.2 (Property tester). Let \mathcal{P} be a property of distributions over $[n]$. A *t-query COND testing algorithm* for \mathcal{P} is a randomized algorithm \mathcal{T} which takes as input $n, \varepsilon \in (0, 1]$, as well as access to COND_D . After making at most $t(\varepsilon, n)$ calls to the oracle, \mathcal{T} either outputs **ACCEPT** or **REJECT**, such that the following holds:

- if $D \in \mathcal{P}$, \mathcal{T} outputs **ACCEPT** with probability at least $2/3$;

²Recall that a non-adaptive tester is an algorithm whose queries do not depend on the answers obtained from previous ones, but only on its internal randomness. Equivalently, it is a tester that can commit “upfront” to all the queries it will make to the oracle.

- if $d_{\text{TV}}(D, \mathcal{P}) \geq \varepsilon$, \mathcal{T} outputs REJECT with probability at least $2/3$.

We observe that the above definitions can be straightforwardly extended to the more general setting of *pairs* of distributions, where given independent access to two oracles COND_{D_1} , COND_{D_2} the goal is to test whether (D_1, D_2) satisfies a property (now a set of pairs of distributions). This will be the case in [Section 3](#), where we will consider equivalence testing, that is the property $\mathcal{P}_{\text{eq}} = \{ (D_1, D_2) : D_1 = D_2 \}$.

2.2 Adaptive Core Testers

In order to deal with adaptivity in our lower bounds, we will use ideas introduced by Chakraborty et al. [[CFG13](#)]. These ideas, for the case of *label-invariant* properties³ allow one to narrow down the range of possible testers and focus on a restricted class of such algorithms called *adaptive core testers*. These core testers do not have access to the full information of the samples they draw, but instead only get to see the relations (inclusions, equalities) between the queries they make and the samples they get. Yet, Chakraborty et al. [[CFG13](#)] show that any tester for a label-invariant property can be converted into a core tester with same query complexity; thus, it is enough to prove lower bounds against this – seemingly – weaker class of algorithms.

We here rephrase the definitions of a core tester and the view they have of the interaction with the oracle (the *configuration* of the samples), tailored to our setting.

Definition 2.3 (Atoms and partitions). Given a family $\mathcal{A} = (A_1, \dots, A_t) \subseteq [n]^t$, the *atoms* generated by \mathcal{A} are the (at most) 2^t distinct sets of the form $\bigcap_{r=1}^t C_r$, where $C_r \in \{A_r, [n] \setminus A_r\}$. The family of all such atoms, denoted $\text{At}(\mathcal{A})$, is the *partition* generated by \mathcal{A} .

This definition essentially captures “all sets (besides the A_i ’s) about which something can be learnt from querying the oracle on the sets of \mathcal{A} .” Now, given such a sequence of queries $\mathcal{A} = (A_1, \dots, A_t)$ and pairs of samples $\mathbf{s} = ((s_1^{(1)}, s_1^{(2)}), \dots, (s_t^{(1)}, s_t^{(2)})) \in A_1^2 \times \dots \times A_t^2$, we would like to summarize “all the label-invariant information available to an algorithm that obtains $((s_1^{(1)}, s_1^{(2)}), \dots, (s_t^{(1)}, s_t^{(2)}))$ upon querying A_1, \dots, A_t for D_1 and D_2 .” This calls for the following definition:

Definition 2.4 (t -configuration). Given $\mathcal{A} = (A_1, \dots, A_t)$ and $\mathbf{s} = ((s_j^{(1)}, s_j^{(2)}))_{1 \leq j \leq t}$ as above, the *t -configuration of \mathbf{s}* consists of the $6t^2$ bits indicating, for all $1 \leq i, j \leq t$, whether

- $s_i^{(k)} = s_j^{(\ell)}$, for $k, \ell \in \{1, 2\}$; and (relations between samples)
- $s_i^{(k)} \in A_j$, for $k \in \{1, 2\}$. (relations between samples and query sets)

In other terms, it summarizes which is the unique atom $S_i \in \text{At}(\mathcal{A})$ that contains $s_i^{(k)}$, and what collisions between samples have been observed.

As aforementioned, the key idea is to argue that, without loss of generality, one can restrict one’s attention to algorithms that only have access to t -configurations, and generate their queries in a specific (albeit adaptive) fashion:

³Recall that a property is label-invariant (or *symmetric*) if it is closed under relabeling of the elements of the support. More precisely, a property of distributions (resp. pairs of distributions) \mathcal{P} is label-invariant if for any distribution $D \in \mathcal{P}$ (resp. $(D_1, D_2) \in \mathcal{P}$) and permutation σ of $[n]$, one has $D \circ \sigma \in \mathcal{P}$ (resp. $(D_1 \circ \sigma, D_2 \circ \sigma) \in \mathcal{P}$).

Definition 2.5 (Core adaptive tester). A *core adaptive distribution tester* for pairs of distributions is an algorithm \mathcal{T} that acts as follows.

- In the i -th phase, based only on its own internal randomness and the configuration of the previous queries A_1, \dots, A_{i-1} and samples obtained $(s_1^{(1)}, s_1^{(2)}), \dots, (s_{i-1}^{(1)}, s_{i-1}^{(2)})$ – whose labels it does not actually know, \mathcal{T} provides:
 - a number k_i^A for each $A \in \text{At}(A_1, \dots, A_{i-1})$, between 0 and $|A \setminus \{s_j^{(1)}, s_j^{(2)}\}_{1 \leq j \leq i-1}|$ (“how many *fresh, not-already-seen* elements of each particular atom A should be included in the next query”)
 - sets $K_i^{(1)}, K_i^{(2)} \subseteq \{1, \dots, i-1\}$ (“which of the samples $s_1^{(k)}, \dots, s_{i-1}^{(k)}$ (whose label is unknown to the tester, but referred to by the index of the query it got them) will be included in the next query”).
- based on these specifications, the next query A_i is drawn (but not revealed to \mathcal{T}) by
 - drawing uniformly at random a set Λ_i in

$$\left\{ \Lambda \subseteq [n] \setminus \{s_j^{(1)}, s_j^{(2)}\}_{1 \leq j \leq i-1} : \forall A \in \text{At}(A_1, \dots, A_{i-1}), |\Lambda \cap A| = k_i^A \right\}.$$

That is, among all sets, containing only “fresh elements,” whose intersection with each atom contains as many elements as \mathcal{T} requires.

- adding the selected previous samples to this set:

$$\Gamma_i \stackrel{\text{def}}{=} \left\{ s_j^{(1)} : j \in K_i^{(1)} \right\} \cup \left\{ s_j^{(2)} : j \in K_i^{(2)} \right\};$$

$$A_i \stackrel{\text{def}}{=} \Lambda_i \cup \Gamma_i.$$

This results in a set A_i , not fully known to \mathcal{T} besides the samples it already got and decided to query again; in which the *labels* of the fresh elements are unknown, but the *proportions* of elements belonging to each atom are known.

- samples $s_i^{(1)} \sim (D_1)_{A_i}$ and $s_i^{(2)} \sim (D_2)_{A_i}$ are drawn (but not disclosed to \mathcal{T}). This defines the i -configuration of A_1, \dots, A_i and $(s_1^{(1)}, s_1^{(2)}), \dots, (s_i^{(1)}, s_i^{(2)})$, which is revealed to \mathcal{T} . Put differently, the algorithm only learns (i) to which of the A_ℓ ’s the new sample belongs, and (ii) if it is one of the previous samples, in which stage(s) and for which of D_1, D_2 it has already seen it.

After $t = t(\varepsilon, n)$ such stages, \mathcal{T} outputs either ACCEPT or REJECT, based only on the configuration of A_1, \dots, A_t and $(s_1^{(1)}, s_1^{(2)}), \dots, (s_t^{(1)}, s_t^{(2)})$ (which is all the information it ever had access to).

Note that in particular, \mathcal{T} does not know the labels of samples it got, nor the actual queries it makes: it knows all about their sizes and sizes of their intersections, but not the actual “identity” of the elements they contain.

2.3 On the use of Yao’s Principle in our lower bounds

We recall Yao’s Principle (e.g., see Chapter 2.2 of [MR95]), a technique which is ubiquitous in the analysis of randomized algorithms. Consider a set S of instances of some problem: what this principle states is that the worst-case expected cost of a randomized algorithm on instances in S

is lower-bounded by the expected cost of the best deterministic algorithm on an instance drawn randomly from S .

As an example, we apply it in a standard way in [Section 4](#): instead of considering a randomized algorithm working on a fixed instance, we instead analyze a *deterministic* algorithm working on a *random* instance. (We note that, importantly, the randomness in the samples returned by the COND oracle is “external” to this argument, and these samples behave identically in an application of Yao’s Principle.)

On the other hand, our application in [Section 3](#) is slightly different, due to our use of adaptive core testers. Once again, we focus on deterministic algorithms working on random instances, and the randomness in the samples is external and therefore unaffected by Yao’s Principle. However, we stress that the randomness in the choice of the set Λ_i is also external to the argument, and therefore unaffected – similar to the randomness in the samples, the algorithm has no control here. Another way of thinking about this randomness is via another step in the distribution over instances: after an instance (which is a pair of distributions) is randomly chosen, we permute the labels on the elements of the distribution’s domain uniformly at random. We note that since the property in question is label-invariant, this does not affect its value. We can then use the model as stated in [Section 2.2](#) for ease of analysis, observing that this can be considered an application of the principle of deferred decisions (as in Chapter 3.5 of [\[MR95\]](#)).

3 A Lower Bound for Equivalence Testing

We prove our main lower bound on the sample complexity of testing equivalence between unknown distributions. We construct two priors \mathcal{Y} and \mathcal{N} over *pairs* of distributions (D_1, D_2) over $[n]$. \mathcal{Y} is a distribution over pairs of distributions of the form (D, D) , namely the case when the distributions are identical. Similarly, \mathcal{N} is a distribution over (D_1, D_2) with $d_{\text{TV}}(D_1, D_2) \geq \frac{1}{4}$. We then show that no algorithm \mathcal{T} making $O(\sqrt{\log \log n})$ queries to $\text{COND}^{D_1}, \text{COND}^{D_2}$ can distinguish between a draw from \mathcal{Y} and \mathcal{N} with constant probability (over the choice of (D_1, D_2) , the randomness in the samples it obtains, and its internal randomness).

We describe the construction of \mathcal{Y} and \mathcal{N} in [Section 3.1](#), and provide a detailed analysis in [Section 3.2](#).

3.1 Construction

We now summarize how a pair of distribution is constructed under \mathcal{Y} and \mathcal{N} . (Each specific step will be described in more detail in the subsequent paragraphs.)

1. Effective Support

- (a) Pick k_b from the set $\{0, 1, \dots, \frac{1}{2} \log n\}$ at random.
- (b) Let $b = 2^{k_b}$ and $m \stackrel{\text{def}}{=} b \cdot n^{1/4}$.

2. Buckets

- (a) ρ and r are chosen with $\sum_{i=1}^{2r} \rho^i = n^{1/4}$.
- (b) Divide $\{1, \dots, m\}$ into intervals B_1, \dots, B_{2r} with $|B_i| = b \cdot \rho^i$.

3. Distributions

- (a) Assign probability mass $\frac{1}{2r}$ uniformly over B_i to generate distribution D_1 .

- (b) (i) Let π_1, \dots, π_r be independent 0/1 with $\Pr(\pi_i = 0) = \frac{1}{2}$.
- (ii) If $\pi_i = 0$, assign probability mass $\frac{1}{4r}$ and $\frac{3}{4r}$ over B_{2i-1} and B_{2i} respectively, else $\frac{3}{4r}$ and $\frac{1}{4r}$ respectively. This generates a distribution D_2 .

4. Support relabeling

- (a) Pick a permutation $\sigma \in S_n$ of the *total* support n .
- (b) Relabel the symbols of D_1 and D_2 according to σ .

5. **Output:** Generate (D_1, D_1) for \mathcal{Y} , and (D_1, D_2) otherwise.

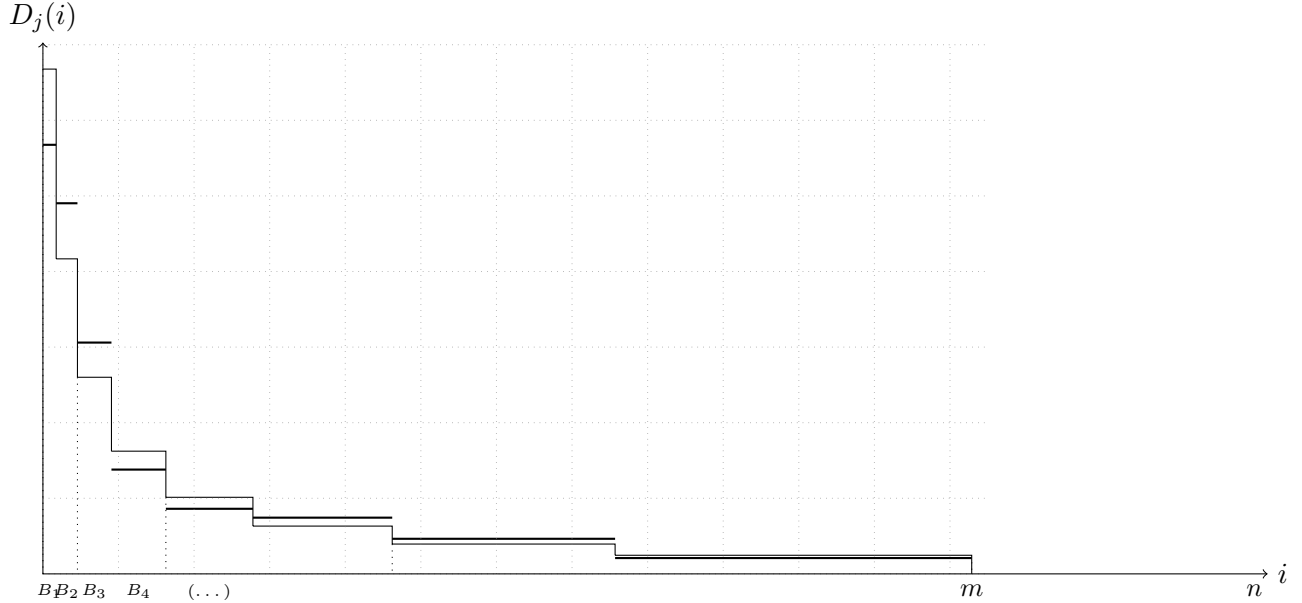


Figure 1: A no-instance (D_1, D_2) (before permutation).

We now describe the various steps of the construction in greater detail.

Effective support. Both D_1 and D_2 , albeit distributions on $[n]$, will have (common) *sparse* support. The support size is taken to be $m \stackrel{\text{def}}{=} b \cdot n^{1/4}$. Note that, from the above definition, m is chosen uniformly at random from products of $n^{1/4}$ with powers of 2, resulting in values in $[n^{1/4}, n^{3/4}]$.

In this step b will act as a random scaling factor. The objective of this random scaling is to induce uncertainty in the algorithm’s knowledge of the true support size of the distributions, and to prevent it from leveraging this information to test equivalence. In fact one can verify that the class of distributions induced for a single value of b , namely all distributions have the same value of m , then one can distinguish the \mathcal{Y} and \mathcal{N} cases with only $O(1)$ conditional queries.

Buckets. Our construction is inspired by the lower bound of Canonne, Ron, and Servedio [CRS15, Theorem 8] for the more restrictive PAIRCOND access model. We partition the support in $2r$ consecutive intervals (henceforth referred to as *buckets*) B_1, \dots, B_{2r} , where the size of the i -th

bucket is $b\rho^i$. We note that r and ρ will be chosen such that $\sum_{i=1}^{2r} b\rho^i = bn^{1/4}$, i.e., the buckets fill the effective support.

Distributions. We output a pair of distributions (D_1, D_2) . Each distribution that we construct is uniform within any particular bucket B_i . In particular, the first distribution assigns the same mass $\frac{1}{2r}$ to each bucket. Therefore, points within B_i have the same probability mass $\frac{1}{(2rb\rho^i)}$. For the \mathcal{Y} case, the second distribution is identical to the first. For the \mathcal{N} case, we pair buckets in r consecutive *bucket-pairs* Π_1, \dots, Π_r , with $\Pi_i = B_{2i-1} \cup B_{2i}$. For the second distribution D_2 , we consider the same buckets as D_1 , but repartition the mass $1/r$ *within* each Π_i . More precisely, in each pair, one of the buckets gets now total probability mass $\frac{1}{4r}$ while the other gets $\frac{3}{4r}$ (so that the probability of every point is either decreased by a factor $\frac{1}{2}$ or increased by $\frac{3}{2}$). The choice of which goes up and which goes down is done uniformly and independently at random for each bucket-pair determined by the random choices of π_i 's.

Random relabeling. The final step of the construction randomly relabels the symbols, namely is a random injective map from $[m]$ to $[n]$. This is done to ensure that no information about the individual symbol labels can be used by the algorithm for testing. For example, without this the algorithm can consider a few symbols from the first bucket and distinguish the \mathcal{Y} and \mathcal{N} cases. As mentioned in [Section 2.3](#), for ease of analysis, the randomness in the choice of the permutation is, in some sense, deferred to the randomness in the choice of Λ_i during the algorithm's execution.

Summary. A no-instance (D_1, D_2) is thus defined by the following parameters: the support size m , the vector $(\pi_1, \dots, \pi_m) \in \{0, 1\}^r$ (which only impacts D_2), and the final permutation σ of the domain. A yes-instance (D_1, D_1) follows an identical process, however, π has no influence on the final outcome. See [Figure 1](#) for an illustration of such a (D_1, D_2) when σ is the identity permutation and thus the distribution is supported over the first m natural numbers.

Values for ρ and r . By setting $r = \frac{\log n}{8 \log \rho} + O(1)$, we have as desired $\sum_{i=1}^{2r} |B_i| = m$ and there is a factor $(1+o(1))n^{1/4}$ between the height of the first bucket B_1 and the one of the last, B_{2r} . It remains to choose the parameter ρ itself; we shall take it to be $2^{\sqrt{\log n}}$, resulting in $r = \frac{1}{8}\sqrt{\log n} + O(1)$. (Note that for the sake of the exposition, we ignore technical details such as the rounding of parameters, e.g. bucket sizes; these can be easily taken care of at the price of cumbersome case analyses, and do not bring much to the argument.)

3.2 Analysis

We now prove our main lower bound, by analyzing the behavior of core adaptive testers (as per [Definition 2.5](#)) on the families \mathcal{Y} and \mathcal{N} from the previous section. In [Section 3.2.1](#), we argue that, with high probability, the sizes of the queries performed by the algorithm satisfy some specific properties. Conditioned upon this event, in [Section 3.2.2](#), we show that the algorithm will get similar information from each query, whether it is running on a yes-instance or a no-instance.

Before moving to the heart of the argument, we state the following fact, straightforward from the construction of our no-instances:

Fact 3.1. *For any (D_1, D_2) drawn from \mathcal{N} , one has $d_{TV}(D_1, D_2) = 1/4$.*

Moreover, as allowing more queries can only increase the probability of success, we hereafter focus on a core adaptive tester that performs exactly $q = \frac{1}{10}\sqrt{\log \log n}$ (adaptive) queries; and will show that it can only distinguish between yes- and no-instances with probability $o(1)$.

3.2.1 Banning “bad queries”

As mentioned in [Section 3.1](#), the draw of a yes- or no-instance involves a random scaling of the size of the support of the distributions, meant to “blind” the testing algorithm. Recall that a testing algorithm is specified by a decision tree, which at step i , specifies how many unseen elements from each atom to include in the query ($\{k_i^A\}$) and which previously seen elements to include in the query (sets $K_i^{(1)}, K_i^{(2)}$, as defined in [Section 2.2](#)), where the algorithm’s choice depends on the observed configuration at that time. Note that, using Yao’s Principle (as discussed in [Section 2.3](#)), these choices are deterministic for a given configuration – in particular, we can think of all $\{k_i^A\}$ and $K_i^{(1)}, K_i^{(2)}$ in the decision tree as being fixed. In this section, we show that all k_i^A values satisfy with high probability some particular conditions with respect to the choice of distribution, where the randomness is over the choice of the support size.

First, we recall an observation from [\[CFG13\]](#), though we modify it slightly to apply to configurations on pairs of distributions and we apply a slightly tighter analysis. This essentially limits the number of states an algorithm could be in by a function of how many queries it makes.

Proposition 3.2. *The number of nodes in a decision tree corresponding to a q -sample algorithm is at most 2^{6q^2+1} .*

Proof. As mentioned in [Definition 2.4](#), an i -configuration can be described using $6i^2$ bits, resulting in at most 2^{6i^2} i -configurations. Since each i -configuration leads us to some node on the i -th level of the decision tree, the total number of nodes can be upper bounded by summing over the number of i -configurations for i ranging from 0 to q , giving us the desired bound. \square

For the sake of the argument, we will introduce a few notions applying to the *sizes* of query sets: namely, the notions of a number being *small*, *large*, or *stable*, and of a vector being *incomparable*. Roughly speaking, a number is small if a uniformly random set of this size does not, in expectation, hit the largest bucket B_{2r} . On the other hand, it is large if we expect such a set to intersect many bucket-pairs (i.e., a significant fraction of the support). The definition of stable numbers is slightly more quantitative: a number β is stable if a random set of size β , for each bucket B_i , either completely misses B_i or intersects it in a number of points very close to the expected number (in this case, we say the set *concentrates* over B_i). Finally, a vector of values (β_j) is incomparable if the union of random sets S_1, \dots, S_m of sizes β_1, \dots, β_m contains (with high probability) an amount of mass $D\left(\bigcup_j S_j\right)$ which is either much smaller or much larger than the probability $D(s)$ of any single element s .

We formalize these concepts in the definitions below. To motivate them, it will be useful to bear in mind that, from the construction described in [Section 3.1](#), the expected intersection of a uniform random set of size β with a bucket B_i is of size $\beta b \rho^i / n$; while the expected probability mass from B_i it contains (under either D_1 or D_2) is $\beta / (2rn)$.

Definition 3.3. Let q be an integer, and let $\varphi = \Theta(q^{5/2})$. A number β is said to be *small* if $\beta < \frac{n}{b\rho^{2r}}$; it is *large* (with relation to some integer q) if $\beta \geq \frac{n}{b\rho^{2r-2\varphi}}$.

Note that the latter condition equivalently means that, in expectation, a set of large size will intersect at least $\varphi + 1$ bucket-pairs (as it hits an expected $2\varphi + 1$ buckets, since $\beta|B_{2r-2\varphi}|/n \geq 1$). From the above definitions we get that, with high probability, a random set of any fixed size will in expectation either hit many or no buckets:

Proposition 3.4. *A number is either small or large with probability $1 - O\left(\frac{\varphi \log \rho}{\log n}\right)$.*

Proof. A number β is neither large nor small if $\frac{\rho^{2\varphi} n}{\beta \rho^{2r}} \leq b \leq \frac{n}{\beta \rho^{2r}}$. The ratio of the endpoints of the interval is $\rho^{2\varphi}$. Since $b = 2^{k_b}$, this implies that at most $\log \rho^{2\varphi} = 2\varphi \log \rho$ values of k_b could result in a fixed number falling in this range. As there are $\Theta(\log n)$ values for k_b , the proposition follows. \square

The next definition characterizes the sizes of query sets for which the expected intersection with any bucket is either close to 0 (less than $1/\alpha$, for some threshold α), or very big (more than α). (It will be helpful to keep in mind that we will eventually use this definition with $\alpha = \text{poly}(q)$.)

Definition 3.5. A number β is said to be α -stable (for $\alpha \geq 1$) if, for each $j \in [2r]$, $\beta \notin \left[\frac{n}{\alpha \beta \rho^j}, \frac{\alpha n}{\beta \rho^j}\right]$. A vector of numbers is said to be α -stable if all numbers it contains are α -stable.

Proposition 3.6. *A number is α -stable with probability $1 - O\left(\frac{r \log \alpha}{\log n}\right)$.*

Proof. Fix some $j \in [2r]$. A number β does not satisfy the definition of α -stability for this j if $\frac{n}{\alpha \beta \rho^j} \leq b \leq \frac{\alpha n}{\beta \rho^j}$. Since $b = 2^{k_b}$, this implies that at most $\log 2\alpha$ values of k_b could result in a fixed number falling in this range. Noting that there are $\Theta(\log n)$ values for k_b and taking a union bound over all $2r$ values for j , the proposition follows. \square

The following definition characterizes the sizes of query sets which have a probability mass far from the probability mass of any individual element. (For the sake of building intuition, the reader may replace ν in the following by the parameter b of the distribution.)

Definition 3.7. A vector of numbers $(\beta_1, \dots, \beta_\ell)$ is said to be (α, τ) -incomparable with respect to ν (for $\tau \geq 1$) if the two following conditions hold.

- $(\beta_1, \dots, \beta_\ell)$ is α -stable.
- Let Δ_j be the minimum $\Delta \in \{0, \dots, 2r\}$ such that $\frac{\beta_j \nu \rho^{2r-\Delta}}{n} \leq \frac{1}{\alpha}$, or $2r$ if no such Δ exists. For all $i \in [2r]$, $\frac{1}{2rn} \sum_{j=1}^{\ell} \beta_j \Delta_j \notin \left[\frac{1}{\tau 2r \nu \rho^i}, \frac{\tau}{2r \nu \rho^i}\right]$.

Recall from the definition of α -stability of a number that a random set of this size either has essentially no intersection with a bucket or “concentrates over it” (i.e., with high probability, the probability mass contained in the intersection with this bucket is very close to the expected value). The above definition roughly captures the following. For any j , Δ_j is the number of buckets that will concentrate over a random set of size β_j . The last condition asks that the total probability mass from D_1 (or D_2) enclosed in the union of m random sets of size $\beta_1, \dots, \beta_\ell$ be a multiplicative factor of τ from the individual probability weight $\frac{1}{2rb\rho^i}$ of a single element from any of the $2r$ buckets.

Proposition 3.8. *Given that a vector of numbers of length ℓ is α -stable, it is (α, q^2) -incomparable with respect to b with probability at least $1 - O\left(\frac{r \log q}{\log n}\right)$.*

Proof. Fix any vector $(\beta_1, \dots, \beta_\ell)$. By the definition above, for each value b such that $(\beta_1, \dots, \beta_\ell)$ is α -stable, we have

$$\beta_j \cdot \frac{\alpha \rho^{2r}}{n} \leq \frac{\rho^{\Delta_j}}{b} < \beta_j \cdot \frac{\alpha \rho^{2r+1}}{n}, \quad j \in [\ell]$$

or, equivalently,

$$\frac{\log \frac{\alpha \beta_j}{n}}{\log \rho} + 2r + \frac{\log b}{\log \rho} \leq \Delta_j < \frac{\log \frac{\alpha \beta_j}{n}}{\log \rho} + 2r + \frac{\log b}{\log \rho} + 1, \quad j \in [\ell].$$

Writing $\lambda_j \stackrel{\text{def}}{=} \frac{\log \frac{\alpha \beta_j}{n}}{\log \rho} + 2r$ for $j \in [\ell]$, we obtain that

$$\sum_{j=1}^{\ell} \beta_j \Delta_j b = b \sum_{j=1}^{\ell} \beta_j (\lambda_j + O(1)) + \frac{b \log b}{\log \rho} \sum_{j=1}^{\ell} \beta_j.$$

- If it is the case that $\log \rho \cdot \sum_{j=1}^{\ell} \beta_j (\lambda_j + O(1)) \ll \log b \cdot \sum_{j=1}^{\ell} \beta_j$. Then, for any fixed $i \in [2r]$, to meet the second item of the definition of incomparability we need $\sum_{j=1}^{\ell} \beta_j \Delta_j b \notin [n/(200q\rho^i), 200qn/\rho^i]$. This is essentially, with the assumption above, requiring that

$$b \log b \notin \left[\frac{n \log \rho}{2q^2 \rho^i \sum_{j=1}^{\ell} \beta_j}, \frac{2q^2 n \log \rho}{\rho^i \sum_{j=1}^{\ell} \beta_j} \right].$$

Recalling that $b \log b = k_b 2^{k_b}$, this means that $O(\log q / \log \log q)$ values of k_b are to be ruled out. (Observe that this is the number of possible ‘‘bad values’’ for b without the condition from the case distinction above; since we add an extra constraint on b , there are at most this many values to avoid.)

- Conversely, if $\log \rho \cdot \sum_{j=1}^{\ell} \beta_j (\lambda_j + O(1)) \gg \log b \cdot \sum_{j=1}^{\ell} \beta_j$ the requirement becomes

$$b \notin \left[\frac{n \log \rho}{2q^2 \rho^i \sum_{j=1}^{\ell} \beta_j (\lambda_j + O(1))}, \frac{2q^2 n \log \rho}{\rho^i \sum_{j=1}^{\ell} \beta_j (\lambda_j + O(1))} \right].$$

ruling out this time $O(\log q)$ values for k_b .

- Finally, the two terms are comparable only if $\log b = \Theta\left(\log \rho \cdot \sum_{j=1}^{\ell} \beta_j (\lambda_j + O(1)) \cdot \left(\sum_{j=1}^{\ell} \beta_j\right)^{-1}\right)$; given that $\log b = k_b$, this rules out this time $O(1)$ values for k_b .

A union bound over the $2r$ possible values of i , and the fact that k_b can take $\Theta(\log n)$ values, complete the proof. \square

We put these together to obtain the following lemma:

Lemma 3.9. *With probability at least $1 - O\left(\frac{2^{6q^2+q}(r \log \alpha + \varphi \log \rho) + 2^{6q^2}(r \log q)}{\log n}\right)$, the following holds for the decision tree corresponding to a q -query algorithm:*

- the size of each atom is α -stable and either large or small;
- the size of each atom, after excluding elements we have previously observed,⁴ is α -stable and either large or small;

⁴More precisely, we mean to say that for each $i \leq q$, for every atom A defined by the partition of (A_1, \dots, A_i) , the values k_i^A and $|A \setminus \{s_1^{(1)}, s_1^{(2)}, \dots, s_{i-1}^{(1)}, s_{i-1}^{(2)}\}| - k_i^A$ are α -stable and either large or small;

- for each i , the vector $(k_i^A)_{A \in \text{At}(A_1, \dots, A_i)}$ is (α, q^2) -incomparable (with respect to b).

Proof. From [Proposition 3.2](#), there are at most 2^{6q^2+1} tree nodes, each of which contains one vector $(k_i^A)_A$, and at most 2^q atom sizes. The first point follows from [Propositions 3.4](#) and [3.6](#) and applying the union bound over all $2^{6q^2+1} \cdot 2 \cdot 2^q$ sizes, where we note the additional factor of 2 comes from either including or excluding the old elements. The latter point follows from [Proposition 3.8](#) and applying the union bound over all 2^{6q^2+1} (k_i^A) vectors. \square

3.2.2 Key lemma: bounding the variation distance between decision trees

In this section, we prove a key lemma on the variation distance between the distribution on leaves of any decision tree, when given access to either an instance from \mathcal{Y} or \mathcal{N} . This lemma will in turn directly yield [Theorem 1.1](#). Hereafter, we set the parameters α (the threshold for stability), φ (the parameter for smallness and largeness) and γ (an accuracy parameter for how well things concentrate over their expected value) as follows:⁵ $\alpha \stackrel{\text{def}}{=} q^7$, $\varphi \stackrel{\text{def}}{=} q^{5/2}$ and $\gamma \stackrel{\text{def}}{=} 1/\varphi = q^{-5/2}$. (Recall further that $q = \frac{1}{10}\sqrt{\log \log n}$.)

Lemma 3.10. *Conditioned on the events of [Lemma 3.9](#), consider the distribution over leaves of any decision tree corresponding to a q -query adaptive algorithm when the algorithm is given a yes-instance, and when it is given a no-instance. These two distributions have total variation distance $o(1)$.*

Proof. This proof is by induction. We will have three inductive hypotheses, $\mathbf{E}_1(t)$, $\mathbf{E}_2(t)$, and $\mathbf{E}_3(t)$. Assuming all three hold for all $t < i$, we prove $\mathbf{E}_1(i)$. Additionally assuming $\mathbf{E}_1(i)$, we prove $\mathbf{E}_2(i)$ and $\mathbf{E}_3(i)$.

Roughly, the first inductive hypothesis states that the query sets behave similarly to as if we had picked a random set of that size. It also implies that whether or not we get an element we have seen before is “obvious” based on past observances and the size of the query we perform. The second states that we never observe two distinct elements from the same bucket-pair. The third states that the next sample is distributed similarly in either a yes-instance or a no-instance. Note that this distribution includes both features which our algorithm can observe (i.e., the atom which the sample belongs to and if it collides with a previously seen sample), as well as those which it can not (i.e., which bucket-pair the observed sample belongs to). It is necessary to show the latter, since the bucket-pair a sample belongs to may determine the outcome of future queries.

More precisely, the three inductive hypotheses are as follows:

- $\mathbf{E}_1(i)$: In either a yes-instance or a no-instance, the following occurs: For an atom S in the partition generated by A_1, \dots, A_i , let $S' = S \setminus \{s_1^{(1)}, s_1^{(2)}, \dots, s_{i-1}^{(1)}, s_{i-1}^{(2)}\}$. For every such S' , let $\ell^{S'}$ be the largest index $\ell \in \{0, \dots, 2r\}$ such that $\frac{|S'|b\rho^\ell}{n} \leq \frac{1}{\alpha}$, or 0 if no such ℓ exists. We claim that $\ell^{S'} \in \{0, \dots, 2r - \varphi - 2\} \cup \{2r\}$, and say S' is small if $\ell^{S'} = 2r$ and large otherwise. Additionally:

- for $j \leq \ell^{S'}$, $|S' \cap B_j| = 0$;
- for $j > \ell^{S'}$, $|S' \cap B_j|$ lies in $[1 - i\gamma, 1 + i\gamma] \frac{|S'|b\rho^j}{n}$.

⁵This choice of parameters is not completely arbitrary: combined with the setting of q , r and ρ , they ensure a total bound $o(1)$ on variation distance and probability of “bad events” as well as a (relative) simplicity and symmetry in the relevant quantities.

Furthermore, let p_1 and p_2 be the probability mass contained in Λ_i and Γ_i , respectively. Then $\frac{p_1}{p_1+p_2} \leq O\left(\frac{1}{q^2}\right)$ or $\frac{p_2}{p_1+p_2} \leq O\left(\frac{1}{q^2}\right)$ (that is, either almost all the probability mass comes from elements which we have not yet observed, or almost all of it comes from previously seen ones).

- $\mathbf{E}_2(i)$: No two elements from the set $\{s_1^{(1)}, s_1^{(2)}, \dots, s_i^{(1)}, s_i^{(2)}\}$ belong to the same bucket-pair.
- $\mathbf{E}_3(i)$: Let T_i^{yes} be the random variable representing the atoms and bucket-pairs⁶ containing $(s_i^{(1)}, s_i^{(2)})$, as well as which of the previous samples they intersect with, when the i -th query is performed on a yes-instance, and define T_i^{no} similarly for no-instances. Then $d_{\text{TV}}(T_i^{\text{yes}}, T_i^{\text{no}}) \leq O\left(\frac{1}{q^2} + \frac{1}{\rho} + \gamma + \frac{1}{\varphi}\right) = o(1)$.

We will show that $\mathbf{E}_1(i)$ holds with probability $1 - O\left(2^i \exp\left(-\frac{2\gamma^2\alpha}{3}\right)\right)$ and $\mathbf{E}_2(i)$ holds with probability $1 - O(i/\varphi)$. Let T^{yes} be the random variable representing the q -configuration and the bucket-pairs containing each of the observed samples in a yes-instance, and define T^{no} similarly for a no-instance. We note that this random variable determines which leaf of the decision tree we reach. By a union bound, coupling argument, and triangle inequality, the total variation distance between T^{yes} and T^{no} will be $O\left(2^q \exp\left(-\frac{2\gamma^2\alpha}{3}\right) + \frac{q^2}{\varphi} + \frac{1}{q} + \frac{q}{\rho} + q\gamma + \frac{q}{\varphi}\right) = o(1)$ (from our choice of α, γ, φ), giving the desired result.

We proceed with the inductive proofs of $\mathbf{E}_1(i)$, $\mathbf{E}_2(i)$, and $\mathbf{E}_3(i)$, noting that the base cases hold trivially for all three of these statements. Throughout this proof, recall that Λ_i is the set of unseen support elements which we query, and Γ_i is the set of previously seen support elements which we query.

Lemma 3.11. *Assuming that $\mathbf{E}_1(t), \mathbf{E}_2(t), \mathbf{E}_3(t)$ hold for all $1 \leq t \leq i-1$, then $\mathbf{E}_1(i)$ holds with probability at least $1 - O\left(2^i \exp\left(-\frac{2\gamma^2\alpha}{3}\right)\right) = 1 - 2^{i-\Omega(q^2)}$.*

Proof. We start with the first part of the statement of $\mathbf{E}_1(i)$, prior to ‘‘Furthermore’’; and let S (and the corresponding S') be any atom as in $\mathbf{E}_1(i)$. First, we note that $\ell^{S'} \in \{0, \dots, 2r - \varphi - 2\} \cup \{2r\}$ since we are conditioning on **Lemma 3.9**: $|S'|$ is α -stable and either large or small, which enforces this condition.

Next, suppose S' is contained in some other atom T generated by A_1, \dots, A_{i-1} , and let $T' = T \setminus \{s_1^{(1)}, s_1^{(2)}, \dots, s_{i-1}^{(1)}, s_{i-1}^{(2)}\}$. Since $|S'| \leq |T'|$, this implies that $\ell^{T'} \leq \ell^{S'}$. We argue about $|T' \cap B_j|$ for three regimes of j :

- The first case is $j \leq \ell^{T'}$. By the inductive hypothesis, $|T' \cap B_j| = 0$, so $|S' \cap B_j| = 0$ with probability 1.
- The next case is $\ell^{T'} < j \leq \ell^{S'}$. Recall from the definition of an core adaptive tester that S' will be chosen uniformly at random from all subsets of T' of the appropriate size. By the inductive hypothesis,

$$\frac{|T' \cap B_j|}{|T'|} \in [1 - (i-1)\gamma, 1 + (i-1)\gamma] \frac{b\rho^j}{n},$$

and therefore

$$\mathbb{E}[|S' \cap B_j|] \in [1 - (i-1)\gamma, 1 + (i-1)\gamma] \frac{|S'| b\rho^j}{n}, \quad \text{implying } \mathbb{E}[|S' \cap B_j|] \leq \frac{2}{\alpha \rho^{\ell^{S'} - j}};$$

⁶If a sample $s_i^{(k)}$ does not belong to any bucket (if the corresponding i -th query did not intersect the support), it is marked in T_i^{yes} with a ‘‘dummy label’’ to indicate so.

where the inequality is by the definition of $\ell^{S'}$ and using the fact that $(i-1)\gamma \leq 1$. Using a Chernoff bound for negatively correlated random variables (as in e.g. [DR96]),

$$\begin{aligned} \Pr[|S' \cap B_j| \geq 1] &= \Pr\left[|S' \cap B_j| \geq \left(1 + \frac{1-\mu}{\mu}\right)\mu\right] \\ &\leq \exp\left(-\frac{(1-\mu)^2}{3\mu}\right) \\ &\leq \exp\left(-\frac{1}{12}\alpha\rho^{\ell^{S'}-j}\right), \end{aligned}$$

where the second inequality holds because $\mu \leq \frac{2}{\alpha\rho^{\ell^{S'}-j}}$ and $(1-\mu)^2 \geq \frac{1}{2}$ for n sufficiently large.

- The final case is $j > \ell^{S'}$. As in the previous one,

$$\mathbb{E}[|S' \cap B_j|] \in [1 - (i-1)\gamma, 1 + (i-1)\gamma] \frac{|S'| b \rho^j}{n}, \quad \text{implying } \mathbb{E}[|S' \cap B_j|] \geq \frac{\alpha\rho^{j-\ell^{S'}-1}}{2};$$

where the inequality is by the definition of $\ell^{S'}$, α -stability, and using the fact that $(i-1)\gamma \leq \frac{1}{2}$. Using again a Chernoff bound for negatively correlated random variables,

$$\begin{aligned} \Pr\left[|S' \cap B_j| - \mathbb{E}[|S' \cap B_j|] \geq \gamma \frac{|S'| b \rho^j}{n}\right] &\leq \Pr[|S' \cap B_j| - \mathbb{E}[|S' \cap B_j|] \geq \gamma 2\mathbb{E}[|S' \cap B_j|]] \\ &\leq 2 \exp\left(-\frac{(2\gamma)^2 \mathbb{E}[|S' \cap B_j|]}{3}\right) \\ &\leq 2 \exp\left(-\frac{2}{3}\gamma^2 \alpha\rho^{j-\ell^{S'}-1}\right) \end{aligned}$$

where the first inequality comes from $2(1 - (i-1)\gamma) \geq 1$, the second one is the Chernoff bound, and the third derives from $\mathbb{E}[|S' \cap B_j|] \geq \frac{\alpha\rho^{j-\ell^{S'}-1}}{2}$.

Combining these, the probability that S' does not satisfy the required conditions is at most

$$\sum_{j \leq \ell^{T'}} 0 + \sum_{\ell^{T'} < j \leq \ell^{S'}} \exp\left(-\frac{1}{12}\alpha\rho^{\ell^{S'}-j}\right) + \sum_{j > \ell^{S'}} 2 \exp\left(-\frac{2}{3}\gamma^2 \alpha\rho^{j-\ell^{S'}-1}\right).$$

This probability is maximized when $\ell^{S'} = \ell^{T'} = 0$, in which case it is

$$\sum_{j=1}^{2r} 2 \exp\left(-\frac{2}{3}\gamma^2 \alpha\rho^{j-1}\right) \leq \sum_{j=1}^{\infty} 2 \exp\left(-\frac{2}{3}\gamma^2 \alpha\rho^{j-1}\right) \leq 3 \exp\left(-\frac{2}{3}\gamma^2 \alpha\right).$$

Taking a union bound over at most 2^i sets gives us the desired probability bound.

Finally, we prove the statement following ‘‘Furthermore’’; this will follow from the definition of incomparability.

- First, we focus on Γ_i . Suppose that Γ_i contains at least one element with positive probability mass (if not, the statement trivially holds). Let p'_2 be the probability mass of the heaviest

element in Γ_i . Since our inductive hypothesis implies that Γ_i has no elements in the same bucket pair, the maximum possible value for p_2 is

$$\begin{aligned} p_2 &\leq p'_2 + \frac{3p'_2}{\rho} + \frac{3p'_2}{\rho^3} + \dots \leq p'_2 + \frac{3p'_2}{\rho} \sum_{k=0}^{\infty} \frac{1}{\rho^{2k}} = \left(1 + \frac{3}{\rho} \frac{\rho^2}{\rho^2 - 1}\right) p'_2 \\ &\leq (1 + o(1))p'_2 \end{aligned}$$

Therefore, $p_2 \in [p'_2, (1 + o(1))p'_2]$. Supposing this heaviest element belongs to bucket j , we can say that $p_2 \in [\frac{1}{2}, (1 + o(1))\frac{3}{2}] \frac{1}{2rb\rho^j}$.

- Next, we focus on Λ_i . Consider some atom A , from which we selected k_A elements which have not been previously observed: call the set of these elements A' . In the first part of this proof, we showed that for each bucket B_k , either $|A' \cap B_k| = 0$ or $|A' \cap B_k| \in [1 - i\gamma, 1 + i\gamma] \frac{|A'|b\rho^k}{n}$. In the latter case, noting that $i\gamma \leq \frac{1}{2}$ and that the probability of an individual element in B_k is within $[1, 3] \frac{1}{4rb\rho^k}$, the probability mass contained by $|A' \cap B_k|$ belongs to $[1, 9] \frac{|A'|}{8rn}$. Recalling the definition of Δ_A as stated in [Definition 3.7](#), as shown earlier in this proof, this non-empty intersection happens for exactly Δ_A bins. Therefore, the total probability mass in Λ_i is in the interval $[\frac{1}{4}, \frac{9}{4}] \frac{1}{2rn} \sum_{A \in \text{At}(A_1, \dots, A_i)} k_i^A \Delta_A$.

Recall that we are conditioning on [Lemma 3.9](#) which states that the vector $(k_i^A)_{A \in \text{At}(A_1, \dots, A_i)}$ is (α, q^2) -incomparable with respect to b . Applying this definition to the bounds just obtained on the probability masses in Λ_i and Γ_i gives the desired result. \square

Lemma 3.12. *Assuming that $\mathbf{E}_1(t), \mathbf{E}_2(t), \mathbf{E}_3(t)$ hold for all $1 \leq t \leq i-1$ and additionally $\mathbf{E}_1(i)$, then $\mathbf{E}_2(i)$ holds with probability at least $1 - O\left(\frac{i}{\varphi}\right)$.*

Proof. We focus on $s_i^{(1)}$. If $s_i^{(1)} \in \Gamma_i$, the conclusion is trivial, so suppose $s_i^{(1)} \in \Lambda_i$. From $\mathbf{E}_1(i)$, no small atom intersects any of the buckets, so condition that $s_i^{(1)}$ belongs to some large atom S . Since we want $s_i^{(1)}$ to fall in a distinct bucket-pair from $2(i-1)+1$ other samples, there are at most $2i-1$ bucket-pairs which $s_i^{(1)}$ should not land in. Using $\mathbf{E}_1(i)$, the maximum probability mass contained in the intersection of these bucket-pairs and S is $(1+i\gamma)(2i-1) \frac{|S|}{rn}$. Similarly, additionally using the definition of a large atom, the minimum probability mass contained in S is $(1-i\gamma)\varphi \frac{|S|}{rn}$. Taking the ratio of these two terms gives an upper bound on the probability of breaking this invariant, conditioned on landing in S , as $O(i/\varphi)$, where we note that $\frac{1+i\gamma}{1-i\gamma} = O(1)$. Since the choice of large atom was arbitrary, we can remove the conditioning. Taking the union bound for $s_i^{(1)}$ and $s_i^{(2)}$ gives the result. \square

Lemma 3.13. *Assuming that $\mathbf{E}_1(t), \mathbf{E}_2(t), \mathbf{E}_3(t)$ hold for all $1 \leq t \leq i-1$ and additionally $\mathbf{E}_1(i)$, then $\mathbf{E}_3(i)$ holds.*

Proof. We focus on some fixed setting of the history of the interaction, i.e. the configuration and the bucket-pairs the past elements belong to, and show that the results of the next query will behave similarly, whether the instance is a **yes**-instance or a **no**-instance. We note that, since we are assuming the inductive hypotheses hold, certain settings which violate these hypotheses are not allowed. We also note that $s_i^{(1)}$ is distributed identically in both instances, so we focus on $s_i \stackrel{\text{def}}{=} s_i^{(2)}$ for the remainder of this proof.

First, we condition that, based on the setting of the past history, s_i will either come from Λ_i or Γ_i – this event happening with probability $1 - O(1/q^2)$.

Proposition 3.14. *In either a yes-instance or a no-instance, s_i will either come from Λ_i or Γ_i with probability $1 - O\left(\frac{1}{q}\right)$, where the choice of which one is deterministic based on the fixed configuration and choice for the bucket-pairs of previously seen elements.*

Proof. This is simply a rephrasing of the portion of $\mathbf{E}_1(i)$ following “Furthermore.” □

By a coupling argument, after conditioning on this event, we must show that the total variation distance in either case is at most $O\left(\frac{1}{\rho} + \gamma + \frac{1}{\varphi}\right) = O\left(\frac{1}{q^{5/2}}\right)$. We break this into two cases, the first being when s comes from Γ_i . In this case, we incur a cost in total variation distance which is $O(1/\rho)$:

Proposition 3.15. *In either a yes-instance or a no-instance, condition that s_i comes from Γ_i . Then, one of the following holds:*

- $|\Gamma_i \cap B_j| = 0$ for all $j \in [2r]$, in which case s_i is distributed uniformly at random from the elements of Γ_i ;
- or $|\Gamma_i \cap B_j| \neq 0$ for some $j \in [2r]$, in which case s_i will be equal to some $s \in \Gamma_i$ with probability $1 - O\left(\frac{1}{\rho}\right)$, where the choice of s is deterministic based on the fixed configuration and choice for the bucket-pairs of previously seen elements.

Proof. The former case follows from the definition of the sampling model. For the latter case, let p be the probability mass of the heaviest element in Γ_i . Since our inductive hypothesis implies that Γ_i has no elements in the same bucket-pair, the maximum possible value for the rest of the elements is

$$\frac{3p}{\rho} + \frac{3p}{\rho^3} + \frac{3p}{\rho^5} + \dots \leq \frac{3p}{\rho} \sum_{k=0}^{\infty} \frac{1}{\rho^{2k}} = \frac{3p}{\rho} \frac{\rho^2}{\rho^2 - 1} = O\left(\frac{p}{\rho}\right).$$

Since the ratio of this value and p is $O\left(\frac{1}{\rho}\right)$, with probability $1 - O\left(\frac{1}{\rho}\right)$ the sample returned is the heaviest element in Γ_i . □

Finally, we examine the case when s comes from Λ_i :

Proposition 3.16. *Condition that s_i comes from Λ_i . Then either:*

- $|\Lambda_i \cap B_j| = 0$ for all $j \in [2r]$, in which case $d_{\text{TV}}(T_i^{\text{yes}}, T_i^{\text{no}}) = 0$;
- or $|\Lambda_i \cap B_j| \neq 0$ for some $j \in [2r]$, in which case $d_{\text{TV}}(T_i^{\text{yes}}, T_i^{\text{no}}) \leq O\left(\gamma + \frac{1}{\varphi}\right) = O\left(\frac{1}{q^{5/2}}\right)$

Proof. The former case follows from the definition of the sampling model – since Λ_i does not intersect any of the buckets, the sample will be labeled as such. Furthermore, the sample returned will be drawn uniformly at random from Λ_i , and the probability of each atom will be proportional to the cardinality of its intersection with Λ_i , in both the yes- and the no-instances.

We next turn to the latter case. Let \mathcal{X} be the event that, if the intersection of Λ_i and some atom A has a non-empty intersection with an odd number of buckets, then s_i does not come from the unpaired bucket. Note that $\mathbf{E}_1(i)$ and the definition of a large atom imply that an unpaired bucket can only occur if the atom intersects at least φ bucket-pairs: conditioned on the sample coming

from a particular atom, the probability that it comes from the unpaired bucket is $O(1/\varphi)$. Since the choice of A was arbitrary, we may remove the conditioning, and note that $\Pr(\mathcal{X}) = 1 - O(1/\varphi)$.

Since

$$\begin{aligned} d_{\text{TV}}(T_i^{\text{yes}}, T_i^{\text{no}}) &\leq d_{\text{TV}}(T_i^{\text{yes}}, T_i^{\text{no}} \mid \mathcal{X}) \Pr(\mathcal{X}) + d_{\text{TV}}(T_i^{\text{yes}}, T_i^{\text{no}} \mid \bar{\mathcal{X}}) \Pr(\bar{\mathcal{X}}) \\ &\leq d_{\text{TV}}(T_i^{\text{yes}}, T_i^{\text{no}} \mid \mathcal{X}) + O(1/\varphi), \end{aligned}$$

it remains to show that $d_{\text{TV}}(T_i^{\text{yes}}, T_i^{\text{no}} \mid \mathcal{X}) \leq O(\gamma)$.

First, we focus on the distribution over atoms, conditioned on \mathcal{X} . Let N^A be the number of bucket-pairs with which A intersects both buckets, i.e., conditioned on \mathcal{X} , the sample could come from $2N^A$ buckets, and let $N \stackrel{\text{def}}{=} \sum_{A \in \text{At}(A_1, \dots, A_i)} N^A$. By $\mathbf{E}_1(i)$, the maximum amount of probability mass that can be assigned to atom A is $\frac{(1+\gamma)|S|N^A/rn}{(1-\gamma)|S|N/rn}$, and the minimum is $\frac{(1-\gamma)|S|N^A/rn}{(1+\gamma)|S|N/rn}$, so the total variation distance in the distribution incurred by this atom is at most $O(\gamma N^A/N)$. Summing over all atoms, we get the desired $O(\gamma)$.

Finally, we bound the distance on the distribution over bucket-pairs, again conditioned on \mathcal{X} . By $\mathbf{E}_1(i)$ only large atoms will contain non-zero probability mass, so condition on the sample coming from some large atom A . Let N^A be the number of bucket-pairs with which A intersects both buckets, i.e., conditioned on \mathcal{X} , the sample could come from $2N^A$ buckets. Using $\mathbf{E}_1(i)$, the maximum amount of probability mass that can be assigned to any intersecting bucket-pair is $\frac{(1+\gamma)\frac{|A|}{rn}}{(1-\gamma)\frac{|A|}{rn}N^A}$, and the minimum is $\frac{(1-\gamma)\frac{|A|}{rn}}{(1+\gamma)\frac{|A|}{rn}N^A}$, so the total variation distance in the distribution incurred by this bucket-pair is at most $O(\gamma \frac{1}{N^A})$. Summing this difference over all N^A bucket-pairs, we get $\frac{2\gamma}{1-\gamma^2} = O(\gamma)$. Since the choice of large atom A was arbitrary, we can remove the conditioning on the choice of atom. The statement follows by applying the union bound on the distribution over bucket-pairs and the distribution over atoms. \square

We note that in both cases, the cost in total variation distance which is incurred is $O(\frac{1}{\rho} + \gamma + \frac{1}{\varphi})$, which implies $\mathbf{E}_3(i)$. \square

This concludes the proof of [Lemma 3.10](#). \square

With this lemma in hand, the proof of the main theorem is straightforward:

Proof of [Theorem 1.1](#). Conditioned on [Lemma 3.9](#), [Lemma 3.10](#) implies that the distribution over the leaves in a **yes**-instance vs. a **no**-instance is $o(1)$. Since an algorithm's choice to accept or reject depends deterministically on which leaf is reached, this bounds the difference between the conditional probability of reaching a leaf which accepts. Since [Lemma 3.9](#) occurs with probability $1 - o(1)$, the difference between the unconditional probabilities is also $o(1)$. \square

4 Lower Bounds for Non-Adaptive Algorithms

In this section, we prove our lower bounds for non-adaptive support size estimation and uniformity testing, rephrased here:

Theorem 1.2 (Non-Adaptive Support Size Estimation). *Any non-adaptive algorithm which, given COND access to an unknown distribution D on $[n]$, estimates the size of its support up to a factor γ must have query complexity $\Omega\left(\frac{\log n}{\log \gamma}\right)$, for any $\gamma \geq \sqrt{2}$.*

Theorem 1.3 (Non-Adaptive Uniformity Testing). *Any non-adaptive algorithm which, given COND access to an unknown distribution D on $[n]$, distinguishes with probability at least $2/3$ between (a) $D = \mathcal{U}$ and (b) $d_{\text{TV}}(D, \mathcal{U}) \geq \frac{1}{4}$, must have query complexity $\Omega(\log n)$.*

These two theorems follow from the same argument, which we outline below before turning to the proof itself.

Structure of the proof. By Yao’s Principle, we consider deterministic tests and study their performance over random distributions, chosen to be uniform over a random subset of carefully picked size. The proof of [Theorem 1.2](#) then proceeds in 3 steps: in [Lemma 4.2](#), we first argue that all queries made by the deterministic tester will (with high probability over the choice of the support size s) behave very “nicely” with regard to s , i.e. not be concentrated around it. Then, we condition on this to bound the total variation distance between the sequence of samples obtained in the two cases we “oppose,” a random distribution from a family \mathcal{D}_1 and the corresponding one from a family \mathcal{D}_2 . In [Lemma 4.3](#) we show that the part of total variation distance due to samples from the small queries is zero, except with probability $o(1)$ over the choice of s . Similarly, [Lemma 4.3](#) allows us to say (comparing both cases to a third “reference” case, a honest-to-goodness uniform distribution over the whole domain, and applying a triangle inequality) that the remaining part of the total variation distance due to samples from the big queries is zero as well, except again with probability $o(1)$. Combining these three lets us to conclude by properties of the total variation distance, as (since the queries are non-adaptive) the distribution over the sequence of samples is a product distribution. (Moreover, applying [Lemma 4.3](#) as a stand-alone enables us, with little additional work, to obtain [Theorem 1.3](#), as our argument in particular implies distributions from \mathcal{D}_1 are hard to distinguish from uniform.)

The families \mathcal{D}_1 and \mathcal{D}_2 . Fix $\gamma \geq \sqrt{2}$ as in [Theorem 1.2](#); writing $\beta \stackrel{\text{def}}{=} \gamma^2$, we define the set

$$\mathcal{S} \stackrel{\text{def}}{=} \left\{ \beta^k n^{1/4} : 0 \leq k \leq \frac{\log n}{2 \log \beta} \right\} = \{n^{1/4}, \beta n^{1/4}, \beta^2 n^{1/4}, \dots, n^{3/4}\}$$

A no-instance $(D_1, D_2) \in \mathcal{D}_1 \times \mathcal{D}_2$ is then obtained by the following process:

- Draw s uniformly at random from \mathcal{S} .
- Pick a uniformly random set $S_1 \subseteq [n]$ of size s , and set D_1 to be uniform on S_1 .
- Pick a uniformly random set $S_2 \subseteq [n]$ of size βs , and set D_2 to be uniform on S_2 .

(Similarly, a yes-instance is obtained by first drawing a no-instance (D_1, D_2) , and discarding D_2 to keep only $(D_1, D_1) \in \mathcal{D}_1 \times \mathcal{D}_1$.)

We will argue that no algorithm can distinguish with high probability between the cases $(D_1, D_2) \sim \mathcal{D}_1 \times \mathcal{D}_2$ and $(D_1, D_2) \sim \mathcal{D}_1 \times \mathcal{D}_1$, by showing that in both cases D_1 and D_2 generates transcripts indistinguishable from those the *uniform* distribution \mathcal{U}_n would yield. This will imply [Theorem 1.2](#), as any algorithm to estimate the support within a multiplicative γ would imply a distinguisher between instances of the form (D_1, D_1) and (D_1, D_2) (indeed, the support sizes of

D_1 and D_2 differ by a factor $\beta = \gamma^2$). (As for [Theorem 1.3](#), observe that any distribution $D_1 \in \mathcal{D}_1$ has constant distance from the uniform distribution on $[n]$, so that a uniformity tester must be able to tell D_1 apart from \mathcal{U}_n .)

Small and big query sets. Let \mathcal{T} be any deterministic non-adaptive algorithm making $q_{\mathcal{T}} \leq q = \frac{1}{100} \frac{\log n}{\log \beta}$ queries. Without loss of generality, we can assume \mathcal{T} makes exactly q queries, and denote them by $A_1, \dots, A_q \subseteq [n]$. Moreover, we let $a_i = |A_i|$, and (again without loss of generality) write $a_1 \geq \dots \geq a_q$.

As a preliminary observation, note that for any $A \subset [n]$ and $0 \leq s \leq n$ we have

$$\mathbb{E}_S |S \cap A| = \frac{|A|s}{n}$$

where the expectation is over a uniform choice of S among all $\binom{n}{s}$ subsets of size s . This observation will lead us to divide the query sets A_i in two groups, depending on the expected size of their intersection with the (random) support.

With this in mind, the following definition will be crucial to our proof. Intuitively, it captures the distribution of *sizes of intersection* of various query sets with the randomly chosen set S .

Definition 4.1. Let $s \geq 1$, and $\mathcal{A} = \{\alpha_1, \dots, \alpha_q\}$ be any set of q integers. For a real number $t > 0$, define

$$C_t(s) \stackrel{\text{def}}{=} \left| \left\{ i \in [q] : \frac{\alpha_i s}{n} \in (\beta^{-t}, \beta^t) \right\} \right|$$

to be the number of *t-hit points* of \mathcal{A} (for s).

The next result will be crucial to prove our lower bounds: it roughly states that if we consider the set of a_i 's and scale them by the random quantity s/n , then the distribution of the random variable generated has an exponential tail with respect to t .

Lemma 4.2 (Hitting Lemma). *Fix \mathcal{A} as in the previous definition. If s is drawn uniformly at random from S , then with probability at least $99/100$.*

$$\sup_{t>0} \frac{C_t(s)}{t} < \frac{2}{100}. \tag{2}$$

We defer the proof of [Lemma 4.2](#) to [Appendix A](#), and proceed to show how to use it in order to bound the contribution of various types of queries to the distinguishability of D_1 and D_2 (in particular, we will apply [Lemma 4.2](#) to the set of query sizes $\{a_1, \dots, a_q\}$.)

Recall that the a_i 's are non-increasing. If $\frac{a_{q'} s}{n} \leq 1$ let $q' \stackrel{\text{def}}{=} q + 1$, otherwise define q' as the largest integer such that $\frac{a_{q'} s}{n} > 1$. We partition the queries made by \mathcal{T} in two: $A_1, \dots, A_{q'}$ are said to be *big*, while $A_{q'+1}, \dots, A_q$ are *small queries*.

Lemma 4.3. *With probability at least $1 - 2^{-10}$, a random distribution from \mathcal{D}_1 or from \mathcal{D}_2 does not intersect with any small query.*

Proof. Let s be the random size drawn for the definition of the instances. We first claim that $\mathbb{E}[|A_{q'+j} \cap S|] \leq \beta^{-50j}$ for all $j \geq 1$, where the expectation is over the uniform choice of set S_1 for

\mathcal{D}_1 . Indeed, by contradiction suppose there is a $j \geq 1$ such that $\mathbb{E}[|A_{q'+j} \cap S|] = \frac{a_{q'+j}s}{n} > \beta^{-50j}$. By definition of q' , for $1 \leq i \leq j$,

$$1 \geq \frac{a_{q'+i}s}{n} > \beta^{-50j}.$$

Therefore, the queries $A_{q'}, A_{q'+1}, \dots, A_{q'+j}$ contribute to C_{50j} , and we obtain $\frac{C_{50j}}{50j} \geq \frac{j}{50j} \geq \frac{2}{100}$, contradicting [Lemma 4.2](#). Thus, the expected intersection can be bounded as follows:

$$\begin{aligned} \mathbb{E}[|(A_{q'+1} \cup A_{q'+2} \cdots \cup A_q) \cap S|] &\leq \mathbb{E}[|A_{q'+1} \cap S|] + \mathbb{E}[|A_{q'+2} \cap S|] + \cdots + \mathbb{E}[|A_q \cap S|] \\ &\leq \beta^{-50} + \beta^{-100} + \dots \\ &\leq 2^{-12}, \end{aligned}$$

for $\beta \geq \sqrt{2}$. From this, we obtain the result holds for \mathcal{D}_1 by Markov's inequality. The same applies to \mathcal{D}_2 with probability of intersection at most 2^{-10} , proving the lemma. \square

We now turn our attention to the sets with *large* intersections. We will show that under \mathcal{D}_1 and \mathcal{D}_2 , the output of querying the sets $A_1, \dots, A_{q'}$ are indistinguishable from simply picking elements uniformly from the sets $A_1, \dots, A_{q'}$. More precisely, we establish the following.

Lemma 4.4. *Let $\eta^* = 2^{-10}$ and $\eta_s = 1/100$; and fix $\ell \in \{1, 2\}$. At least an $1 - \eta_s$ fraction of elements $s_1, \dots, s_{q'} \in A_1 \times A_2, \dots, A_{q'}$ satisfy*

$$\Pr_j[(s_1, \dots, s_{q'})] \in [1 - 5\eta^*, 1 + 5\eta^*] \cdot \frac{1}{|A_1| \dots |A_{q'}|},$$

where $\Pr_j[(s_1, \dots, s_{q'})]$ denotes the probability that $s_1 \dots s_{q'}$ are the output of the queries $A_1, \dots, A_{q'}$ under COND_{D_ℓ} .

As this holds for most distributions in both \mathcal{D}_1 and \mathcal{D}_2 , this implies the queries are indistinguishable from the product distribution over $A_1 \times A_2, \dots, A_{q'}$ (which is the one induced by the same queries on the uniform distribution over $[n]$) in either case, with probability at least $1 - \eta_s - 5\eta^*$.

Proof of Lemma 4.4. From standard Chernoff-like concentration bounds for negatively correlated random variables (see e.g. [\[DR96, Theorem 4.3\]](#)), we obtain

Claim 4.5. *Suppose A is a set of size a , and S is a uniformly chosen random set of size s . Then, for all $\eta \in (0, 1]$, we have $\Pr[|A \cap S| > (1 + \eta)\frac{as}{n}] < e^{-\eta^2 \cdot \frac{as}{3n}}$ and $\Pr[|A \cap S| < (1 - \eta)\frac{as}{n}] < e^{-\eta^2 \cdot \frac{as}{3n}}$.*

We use this to show that indeed all the $|A_i \cap S|$ concentrate around their expected values for $1 \leq j \leq q'$. First note that, as a consequence of [Lemma 4.2](#), it is the case that these expected values satisfy $a_{q'-j}s/n \geq \beta^{50(j+1)}$ for every $0 \leq j \leq q' - 1$ (with probability at least 99/100). Conditioning on this, we first invoke [Claim 4.5](#) on A_j with $\eta = 3 \cdot \beta^{20(j+1)}$, and then apply a union bound to obtain

$$\Pr \left[\exists j \in [q'] \text{ s.t. } |A_j \cap S| \notin \left[1 - 4 \cdot \beta^{-20(j+1)}, 1 + 4 \cdot \beta^{-20(j+1)} \right] \cdot \frac{a_j s}{n} \right] < e^{-\beta^{10}} \quad (3)$$

i.e., with high probability all intersections simultaneously concentrate around their expected values.

Note that since s is at most $n^{3/4}$, each A_i under consideration has size at least $n\beta^{50}/n^{3/4} > n^{1/4}$. Therefore, the probability that a random selection of elements from $A_1, \dots, A_{q'}$ exhibits no collision is at least

$$\prod_{i=1}^{q'} \frac{|A_i| - q'}{|A_i|} \geq \left(1 - \frac{q'}{n^{1/4}}\right)^{q'} \geq 1 - \frac{(q')^2}{n^{1/4}} > 1 - \frac{\log^2 n}{n^{1/4}}.$$

We henceforth condition on this event.

Let $N = \binom{n}{s}$ be the number of outcomes for the set S . We write $N_0 \geq N(1 - e^{-\beta^{10}})$ for the number of such sets for which (3) holds. Let $s_1^{q'}$ denote $s_1 \dots s_{q'}$. For a set of distinct $(s_1, \dots, s_{q'}) \in A_1 \times \dots \times A_{q'}$, let $N(s_1^{q'}) = \binom{n-q'}{s-q'}$ be the number of sets of size s that contain $s_1^{q'}$, and let $N_0(s_1^{q'})$ of them satisfy (3).

By Markov's inequality, with probability at least $1 - e^{-\beta^9}$, for a randomly chosen $s_1^{q'}$ we have $N_0(s_1^{q'})/N(s_1^{q'}) > 1 - e^{-\beta^9}$. For any such $s_1^{q'}$,

$$\begin{aligned} \Pr[s_1^{q'}] &\geq \frac{N_0(s_1^{q'})}{N} \cdot \prod_{i=1}^{q'} \frac{1}{|A_i \cap S|} \geq (1 - e^{-\beta^9}) \frac{s(s-1)\dots(s-q'+1)}{n(n-1)\dots(n-q'+1)} \cdot (1 - 4 \cdot \beta^{-19}) \prod_{i=1}^{q'} \frac{n}{a_i s} \\ &\geq (1 - 6 \cdot \beta^{-19}) \prod_{i=1}^{q'} \frac{1}{a_i}, \end{aligned}$$

for large n and as $|S| > n^{1/4}$. Since the sum of probabilities of elements is at most 1, the other side of the inequality in Lemma 4.4 follows. \square

Proof of Theorem 1.2 and Theorem 1.3. Let T_1 (resp. T_2, T_U) be the distribution over transcripts generated by the queries A_1, \dots, A_q when given conditional access to the distribution D_1 from a no-instance (resp. D_2 , resp. uniform \mathcal{U}_n); that is, a distribution over q -tuples in $A_1 \times \dots \times A_q$. Since the queries were non-adaptive, we can break T_1 (and similarly for T_2, T_U) in $T_1^{\text{big}} \times T_1^{\text{small}}$ according to q' , and use Lemma 4.4 and Lemma 4.3 separately to obtain $d_{\text{TV}}(T_1, T_2) \leq \eta_s + \eta^* + 2^{-10} < 1/50$ and $d_{\text{TV}}(T_1, T_U) \leq \eta_s + \eta^* + 2^{-10} < 1/50$ (for the latter, recalling that queries that do not intersect the support receive samples exactly uniformly distributed in the query set) – thus establishing both theorems. \square

5 An Upper Bound for Support Size Estimation

In this section, we prove our upper bound for constant-factor support size estimation, reproduced below:

Theorem 1.4 (Adaptive Support Size Estimation). *Let $\tau > 0$ be any constant. There exists an adaptive algorithm which, given COND access to an unknown distribution D on $[n]$ which has minimum non-zero probability τ/n and accuracy parameter ε makes $\tilde{O}((\log \log n)/\varepsilon^3)$ queries to the oracle and outputs a value $\tilde{\omega}$ such that the following holds. With probability at least $2/3$, $\tilde{\omega} \in [\frac{1}{1+\varepsilon} \cdot \omega, (1+\varepsilon) \cdot \omega]$, where $\omega = |\text{supp}(D)|$.*

Before describing and analyzing our algorithm, we shall need the following results, that we will use as subroutines: the first one will help us detecting when the support is already dense. The second, assuming the support is sparse enough, will enable us to find an element with zero

probability mass, which can afterwards be used as a “reference” to verify whether any given element is inside or outside the support. Finally, the last one will use such a reference point to check whether a candidate support size σ is smaller or significantly bigger than the actual support size.

Lemma 5.1. *Given $\tau > 0$ and COND access to a distribution D such that each support element has probability at least τ/n , as well as parameters $\varepsilon \in (0, 1/2), \delta \in (0, 1)$, there exists an algorithm TESTSMALLSUPPORT (Algorithm 2) that makes $\tilde{O}(1/(\tau\varepsilon^2) + 1/\tau^2) \cdot \log(1/\delta)$ queries to the oracle, and satisfies the following. (i) If $\text{supp}(D) \geq (1 - \varepsilon/2)n$, then it outputs ACCEPT with probability at least $1 - \delta$; (ii) if $\text{supp}(D) \leq (1 - \varepsilon)n$, then it outputs REJECT with probability at least $1 - \delta$.*

Lemma 5.2. *Given COND access to a distribution D , an upper bound $m < n$ on $\text{supp}(D)$, as well as parameter $\delta \in (0, 1)$, there exists an algorithm GETNONSUPPORT (Algorithm 3) that makes $\tilde{O}(\log^2 \frac{1}{\delta} \log^{-2} \frac{n}{m})$ queries to the oracle, and returns an element $r \in [n]$ such that $r \notin \text{supp}(D)$ with probability at least $1 - \delta$.*

Lemma 5.3. *Given COND access to a distribution D , inputs $\sigma \geq 2$ and $r \notin \text{supp}(D)$, as well as parameters $\varepsilon \in (0, 1/2), \delta \in (0, 1)$, there exists an algorithm ISATMOSTSUPPORTSIZE (Algorithm 4) that makes $\tilde{O}(1/\varepsilon^2) \log(1/\delta)$ queries to the oracle, and satisfies the following. The algorithm returns either yes or no, and (i) if $\sigma \geq \text{supp}(D)$, then it outputs yes with probability at least $1 - \delta$; (ii) if $\sigma > (1 + \varepsilon)\text{supp}(D)$, then it outputs no with probability at least $1 - \delta$.*

We defer the proofs of these 3 lemmas to the next subsections, and now turn to the proof of the theorem.

Proof. The algorithm is given in Algorithm 1, and at a high-level works as follows: it first checks whether the support size is big (an $1 - O(\varepsilon)$ fraction of the domain), in which case it can already stop and return a good estimate. If this is not the case, however, then the support is sparse enough to efficiently find an element r outside the support, by taking a few uniform points, comparing and ordering them by probability mass (and keeping the lightest). This element r can then be used as a reference point in a (doubly exponential) search for a good estimate: for each guess $\tilde{\omega}$, a random subset S of size roughly $\tilde{\omega}$ is taken, a point x is drawn from D_S , and x is compared to r to check if $D(x) > 0$. If so, then S intersects the support, meaning that $\tilde{\omega}$ is an upper bound on ω ; repeating until this is no longer the case results in an accurate estimate of ω .

Algorithm 1 ESTIMATESUPPORT $_D$

```

1: if TESTSMALLSUPPORT $_D(\varepsilon, \frac{1}{10})$  returns ACCEPT then return  $\tilde{\omega} \leftarrow (1 - \varepsilon^2)n$ 
2: end if
3: Call GETNONSUPPORT $_D((1 - \frac{\varepsilon}{2})n, \frac{1}{10})$  to obtain a non-support reference point  $r$ .
4: for  $j$  from 0 to  $\log_{1+\varepsilon} \log_{1+\varepsilon} n$  do
5:   Set  $\tilde{\omega} \leftarrow (1 + \varepsilon)^{(1+\varepsilon)^j}$ .
6:   Call ISATMOSTSUPPORTSIZE $_D(\tilde{\omega}, r, \varepsilon, \frac{1}{100 \cdot (j+1)^2})$  to check if  $\tilde{\omega}$  is an upper bound on  $\omega$ .
7:   if the call returned no then
8:     Perform a binary search on  $\{(1 + \varepsilon)^{j-1}, \dots, (1 + \varepsilon)^j\}$  to find  $i^*$ , the smallest  $i \geq 2$  such
     that ISATMOSTSUPPORTSIZE $_D((1 + \varepsilon)^i, r, \varepsilon, \frac{1}{10(j+1)})$  returns no.
9:     return  $\tilde{\omega} \leftarrow (1 + \varepsilon)^{i^*-1}$ .
10:  end if
11: end for

```

In the rest of this section, we formalize and rigorously argue the above. Conditioning on each of the calls to the subroutines TESTSMALLSUPPORT, GETNONSUPPORT and ISATMOSTSUPPORTSIZE being correct (which overall happens except with probability at most $1/10+1/10+\sum_{j=1}^{\infty} 1/(100j^2)+1/10 < 1/3$ by a union bound), we show that the output $\tilde{\omega}$ of ESTIMATESUPPORT is indeed within a factor $(1 + \varepsilon)$ of ω .

- If the test on Step 1 passes, then by Lemma 5.1 we must have $\text{supp}(D) > (1 - \varepsilon)n$. Thus, the estimate we output is correct, as $[(1 - \varepsilon)n, n] \subseteq [\tilde{\omega}/(1 + \varepsilon), (1 + \varepsilon)\tilde{\omega}]$.
- Otherwise, if it does not then by Lemma 5.1 it must be the case that $\text{supp}(D) < (1 - \varepsilon/2)n$.

Therefore, if we reach Step 3 then $(1 - \varepsilon/2)n$ is indeed an upper bound on ω , and GETNONSUPPORT will return a point $r \notin \text{supp}(D)$ as expected. The analysis of the rest of the algorithm is straightforward: from the guarantee of ISATMOSTSUPPORTSIZE, the binary search will be performed for the first index j such that $\omega \in [(1 + \varepsilon)^{(1+\varepsilon)^{j-1}}, (1 + \varepsilon)^{(1+\varepsilon)^j}]$; and will be on a set of $(1 + \varepsilon)^{j-1}$ values. Similarly, for the value i^* eventually obtained, it must be the case that $(1 + \varepsilon)^{i^*} > \omega$ (by contrapositive, as no was returned by the subroutine) but $(1 + \varepsilon)^{i^*-1} \leq (1 + \varepsilon)\omega$ (again, as the subroutine returned yes). But then, $\tilde{\omega} = (1 + \varepsilon)^{i^*-1} \in (\omega/(1 + \varepsilon), (1 + \varepsilon)\omega]$ as claimed.

Query complexity. The query complexity of our algorithm originates from the following different steps:

- the call to TESTSMALLSUPPORT, which from Lemma 5.1 costs $\tilde{O}(1/\varepsilon^2)$ queries;
- the call to GETNONSUPPORT, on Step 3, that from the choice of the upper bound also costs $\tilde{O}(1/\varepsilon^2)$ queries;
- the (at most) $\log_{1+\varepsilon} \log_{1+\varepsilon} n = O((\log \log n)/\varepsilon)$ calls to ISATMOSTSUPPORTSIZE on Step 6. Observing that the query complexity of ISATMOSTSUPPORTSIZE is only $\tilde{O}(1/\varepsilon^2) \cdot \log(1/\delta)$, and from the choice of $\delta = \frac{1}{(j+1)^2}$ at the j -th iteration this step costs at most

$$\tilde{O}\left(\frac{1}{\varepsilon^2}\right) \cdot \sum_{j=1}^{\log_{1+\varepsilon} \log_{1+\varepsilon} n} O(\log(j^2)) = \tilde{O}\left(\frac{1}{\varepsilon^2} \log_{1+\varepsilon} \log_{1+\varepsilon} n\right) = \tilde{O}\left(\frac{1}{\varepsilon^3} \log_{1+\varepsilon} \log_{1+\varepsilon} n\right)$$

queries.

- Similarly, Step 8 results in at most $j \leq \log \log n$ calls to ISATMOSTSUPPORTSIZE with δ set to $1/(10(j + 1))$, again costing $\tilde{O}\left(\frac{1}{\varepsilon^2}\right) \cdot \log j = \tilde{O}\left(\frac{1}{\varepsilon^2} \log_{1+\varepsilon} \log_{1+\varepsilon} n\right) = \tilde{O}\left(\frac{1}{\varepsilon^3} \log \log n\right)$ queries.

Gathering all terms, the overall query complexity is $\tilde{O}\left(\frac{\log \log n}{\varepsilon^3}\right)$, as claimed. \square

5.1 Proof of Lemma 5.1.

Hereafter, we assume without loss of generality that $\tau < 2$: indeed, if $\tau \geq 2$ then the support is of size at most $n/2$, and it suffices to output REJECT to meet the requirements of the lemma. We will rely on the (easy) fact below, which ensures that any distribution with dense support and minimum non-zero probability τ/n put significant mass on “light” elements:

Fact 5.4. *Fix any $\varepsilon \in [0, 1)$. Assume D satisfies both $\text{supp}(D) \geq (1 - \varepsilon)n$ and $D(x) \geq \tau/n$ for $x \in \text{supp}(D)$. Then, setting $L_\varepsilon \stackrel{\text{def}}{=} \{x \in [n] : D(x) \in [\tau/n, 2/n]\}$, we have $|L_\varepsilon| \geq (1/2 - \varepsilon)n$ and $D(L_\varepsilon) \geq (1/2 - \varepsilon)\tau$.*

Proof. As the second claim follows directly from the first and the minimum mass of elements of L_ε , it suffices to prove that $|L_\varepsilon| \geq (1/2 - \varepsilon)n$. This follows from observing that

$$1 = D([n]) \geq D([n] \setminus L_\varepsilon) \geq (|\text{supp}(D)| - |L_\varepsilon|) \frac{2}{n} \geq 2(1 - \varepsilon) - \frac{2|L_\varepsilon|}{n}$$

and rearranging the terms. \square

Description and intuition. The algorithm (as described in [Algorithm 2](#)) works as follows: it first takes enough uniformly distributed samples s_1, \dots, s_m to get (with high probability) an accurate enough fraction of them falling in the support to distinguish between the two cases. The issue is now to detect those m_j 's which indeed are support elements; note that we do not care about underestimating this fraction in case (b) (when the support is at most $(1 - \varepsilon)n$, but importantly do not want to underestimate it in case (a) (when the support size is at least $(1 - \varepsilon/2)n$). To perform this detection, we take constantly many samples *according to* D (which are therefore ensured to be in the support), and use pairwise conditional queries to sort them by increasing probability mass (up to approximation imprecision), and keep only the lightest of them, t . In case (a), we now from [Fact 5.4](#) that with high probability our t has mass in $[1/n, 2/n]$, and will therefore be either much lighter than or comparable to *any* support element: this will ensure that in case (a) we do detect all of the m_j 's that are in the support.

This also works in case (b), even though [Fact 5.4](#) does not give us any guarantee on the mass of t . Indeed, either t turns out to be light (and then the same argument ensures us our estimate of the number of “support” m_j 's is good), or t is too heavy – and then our estimate will end up being smaller than the true value. But this is fine, as the latter this only means we will reject the distribution (as we should, since we are in the small-support case).

Correctness. Let η be the fraction of the s_j 's that are in the support of the distribution. By a multiplicative Chernoff bound and a suitable constant in our choice of m , we get that (i) if $\text{supp}(D) \geq 1 - \varepsilon/2$, then $\Pr[\eta < 1 - 3\varepsilon/4] \leq 1/12$, while (ii) if $\text{supp}(D) \leq 1 - \varepsilon/2$, then $\Pr[\eta \geq 1 - 3\varepsilon/4] \leq 1/12$. We hereafter condition on this (i.e., η being a good enough estimate). We also condition on all calls to COMPARE yielding results as per specified, which by a union bound overall happens except with probability $1/12 + 1/12 = 1/6$, and break the rest of the analysis in two cases.

- (a) Since the support size ω is in this case at least $(1 - \varepsilon/2)n$, from [Fact 5.4](#) we get that $D(L_{\varepsilon/2}) \geq \frac{1 - \varepsilon}{2} \tau \geq \frac{\tau}{4}$. Therefore, except with probability at most $(1 - \tau/4)^k < 1/12$, at least one of the t_j 's will belong to $L_{\varepsilon/2}$. When this happens, and by the choice of parameters in the calls to COMPARE, we get that $t \in L_{\varepsilon/2}$; that is $D(t) \in [\tau/n, 2/n]$. But then the calls to the routine on [Step 15](#) will always return either a value (since t is “comparable” to all $x \in L_{\varepsilon/2}$ – i.e., has probability within a factor $2/\tau$ of them) or High (possible for those s_j 's that have weight greater than $2/n$), unless s_j has mass 0 (that is, is not in the support). Therefore, the fraction of points marked as support is exactly η , which by the foregoing discussion is at least $1 - 3\varepsilon/4$: the algorithm returns ACCEPT at [Step 20](#).
- (b) Conversely, if $\omega \leq (1 - \varepsilon)n$, there will be a fraction $1 - \eta > 3\varepsilon/4$ of the s_j 's having mass 0. However, no matter what t is it will still be in the support and therefore have $D(t) \geq \tau/n$: for these s_j 's, the call to COMPARE on [Step 15](#) can thus only return Low. This means that

Algorithm 2 TESTSMALLSUPPORT_D

Require: COND access to D ; accuracy parameter $\varepsilon \in (0, 1/2)$, threshold $\tau > 0$, probability of failure δ

```
1: Repeat the following  $O(\log(1/\delta))$  times and output the majority vote.
2: loop
3:   Draw  $m \stackrel{\text{def}}{=} \Theta\left(\frac{1}{\varepsilon^2}\right)$  independent samples  $s_1, \dots, s_m \sim \mathcal{U}$ .
4:   Draw  $k \stackrel{\text{def}}{=} \Theta\left(\frac{1}{\tau}\right)$  independent samples  $t_1, \dots, t_k \sim D$ .
5:   for all  $1 \leq i < j \leq k$  do ▷ Order the  $t_j$ 's
6:     Call COMPARE( $\{t_i\}, \{t_j\}, \eta = \frac{1}{2}, K = 2, \frac{1}{4k^2}$ ) to get a 2-approx.  $\rho$  of  $\frac{D(t_j)}{D(t_i)}$ , High or Low.
7:     if COMPARE returned High or a value  $\rho$  then
8:       Record  $t_i \preceq t_j$ 
9:     else
10:      Record  $t_j \prec t_i$ 
11:    end if
12:  end for
13:  Set  $t$  to be (any of the) smallest  $t_j$ 's, according to  $\preceq$ .
14:  for all  $1 \leq j \leq m$  do ▷ Find the fraction of support elements among the  $m_j$ 's
15:    Call COMPARE( $\{t\}, \{s_j\}, \eta = \frac{1}{2}, K = \frac{2}{\tau}, \frac{1}{4m}$ ) to get either a value  $\rho$ , High or Low.
16:    if COMPARE returned High or a value  $\rho \geq 1/2$  then
17:      Record  $s_j$  as “support.”
18:    end if
19:  end for
20:  if the number of  $s_j$ 's marked “support” is at least  $(1 - \frac{3}{4}\varepsilon)m$  then return ACCEPT
21:  else return REJECT
22:  end if
23: end loop
```

there can only be less than $(1 - \frac{3}{4}\varepsilon)m$ points marked “support” among the s_j 's, and hence that the algorithm will output REJECT as it should.

Overall, the inner loop of the algorithm thus only fails with probability at most $1/12 + 1/6 + 1/12 = 1/3$ (respectively for η failing to be a good estimate, the calls to COMPARE failing to yield results as guaranteed, or no t_j hitting $L_{\varepsilon/2}$ in case (a)). Repeating independently $\log(1/\delta)$ times and taking the majority vote boosts the probability of success to $1 - \delta$.

Query complexity. The sample complexity comes from the k^2 calls on Step 4 (each costing $O(\log k)$ queries) and the m calls on Step 15 (each costing $O(\frac{1}{\tau} \log m)$ queries). By the setting of m and because of the $\log(1/\delta)$ repetitions, this results in an overall query complexity $O\left(\left(\frac{1}{\tau^2} \log \frac{1}{\tau} + \frac{1}{\tau\varepsilon^2} \log \frac{1}{\varepsilon}\right) \log \frac{1}{\delta}\right)$. \square

5.2 Proof of Lemma 5.2.

As described in Algorithm 3, the subroutine is fairly simple: using its knowledge of an upperbound on the support size, it takes enough uniformly distributed samples to have (with high probability)

at least one falling outside the support. Then, it uses the conditional oracle to “order” these samples according to their probability mass, and returns the lightest of them – i.e., one with zero probability mass.

Algorithm 3 GETNONSUPPORT $_D(m, \delta)$

Require: COND access to D ; upper bound m on $\text{supp}(D)$, probability of failure δ

Ensure: Returns $r \in [n]$ such that, with probability at least $1 - \delta$, $r \notin \text{supp}(D)$

```

1: Set  $k \stackrel{\text{def}}{=} \lceil \log \frac{2}{\delta} \log^{-1} \frac{n}{m} \rceil$ .
2: Draw independently  $k$  points  $s_1, \dots, s_k \sim \mathcal{U}_{[n]}$ 
3: for all  $1 \leq i < j \leq k$  do
4:   Call COMPARE( $\{s_i\}, \{s_j\}, \eta = \frac{1}{2}, K = 2, \frac{\delta}{2k^2}$ ) to get a 2-approx.  $\rho$  of  $\frac{D(s_j)}{D(s_i)}$ , High or Low.
5:   if COMPARE returned High or a value  $\rho$  then
6:     Record  $s_i \preceq s_j$ 
7:   else
8:     Record  $s_j \prec s_i$ 
9:   end if
10: end for
11: return  $\arg \min_{\preceq} \{s_1, \dots, s_k\}$  ▷ Return (any) minimal element for  $\preceq$ .

```

Correctness. It is straightforward to see that provided at least one of the s_j 's falls outside the support and that all calls to COMPARE behave as expected, then the procedure returns one of the “lightest” s_j 's, i.e. a non-support element. By a union bound, the latter holds with probability at least $1 - \delta/2$; as for the former, since m is by assumption an upper bound on the support size it holds with probability at least $1 - (m/n)^k \geq 1 - \delta/2$ (from our setting of k). Overall, the procedure's output is correct with probability at least $1 - \delta$, as claimed.

Query complexity. The query complexity of GETNONSUPPORT is due to the $\binom{k}{2}$ calls to COMPARE, and is therefore $O\left(k^2 \log \frac{k}{\delta}\right)$ because of our setting for η and K (which is in turn $\tilde{O}\left(\log^2 \frac{1}{\delta} \log^{-2} \frac{n}{m}\right)$). (In our case, we shall eventually take $m = (1 - \varepsilon/2)n$ and $\delta = 1/10$, thus getting $k = O(1/\varepsilon)$ and a query complexity of $\tilde{O}(1/\varepsilon^2)$.) \square

5.3 Proof of Lemma 5.3

Our final subroutine, described in Algorithm 4, essentially derives from the following observation: a random set S of size (approximately) σ obtained by including independently each element of the domain with probability $1/\sigma$ will intersect the support on ω/σ points on expectation. What we *can* test given our reference point $r \notin \text{supp}(D)$, however, is only whether $S \cap \text{supp}(D) = \emptyset$. But this is enough, as by repeating sufficiently many times (taking a random S and testing whether it intersects the support at all) we can distinguish between the two cases we are interested in. Indeed, the expected fraction of times S includes a support element in either cases is known to the algorithm and differs by roughly $\Omega(\varepsilon)$, so $O(1/\varepsilon^2)$ repetitions are enough to tell the two cases apart.

Algorithm 4 ISATMOSTSUPPORTSIZE $_D(\sigma, r, \varepsilon, \delta)$

Require: COND access to D ; size $\sigma \geq 2$, non-support element r , accuracy ε , probability of failure δ

Ensure: Returns, with probability at least $1 - \delta$, **yes** if $\sigma \leq |\text{supp}(D)|$ and **no** if $\sigma > (1 + \varepsilon)|\text{supp}(D)|$.

- 1: Set $\alpha \leftarrow \left(1 - \frac{1}{\sigma}\right)^\sigma \in [\frac{1}{4}, e^{-1}]$, $\tau \leftarrow \alpha(\alpha^{-\frac{\varepsilon}{2}} - 1) = \Theta(\varepsilon)$.
- 2: Repeat the following $O(\log(1/\delta))$ times and output the majority vote.
- 3: **loop**
- 4: **for** $m = O\left(\frac{1}{\tau^2}\right)$ times **do**
- 5: Draw a subset $S \subseteq [n]$ by including independently each $x \in [n]$ with probability $1/\sigma$.
- 6: Draw $x \sim D_S$.
- 7: Call COMPARE($\{x\}, \{r\}, \eta = \frac{1}{2}, K = 1, \frac{1}{100m}$) \triangleright **Low** if $S \cap \text{supp}(D) \neq \emptyset$; $\rho \in [\frac{1}{2}, 2)$ o.w.
- 8: Record **yes** if COMPARE returned **Low**, **no** otherwise.
- 9: **end for**
- 10: **return yes** if at least $m(\alpha + \frac{\tau}{2})$ “**yes**”’s were recorded, **no** otherwise. \triangleright Thresholding.
- 11: **end loop**

Correctness. We condition on all calls to COMPARE being correct: by a union bound, this overall happens with probability at least 99/100. We shall consider the two cases $\sigma \leq \omega$ and $\sigma > (1 + \varepsilon)\omega$, and focus on the difference of probability p of recording **yes** on Step 8 between the two, in any fixed of the m iterations. In both cases, note p is exactly $(1 - 1/\sigma)^\omega$.

- If $\sigma \leq \omega$, then we have $p \leq \left(1 - \frac{1}{\sigma}\right)^\sigma = \alpha$.
- If $\sigma > (1 + \varepsilon)\omega$, then $p > \left(1 - \frac{1}{\sigma}\right)^{\sigma/(1+\varepsilon)} > \left(1 - \frac{1}{\sigma}\right)^{\sigma(1-\varepsilon/2)} = \alpha^{1-\varepsilon/2}$.

As $\alpha \in [\frac{1}{4}, e^{-1}]$, the difference between the two is $\tau = \alpha(\alpha^{-\varepsilon/2} - 1) = \Theta(\varepsilon)$. Thus, repeating the atomic test of Step 4 $O(1/\tau^2)$ before thresholding at Step 10 yields the right answer with constant probability, then brought to $1 - \delta$ by the outer repeating and majority vote.

Query complexity. Each call to COMPARE at Step 7 costs $O(\log m)$ queries, and is in total repeated $O(m \log(1/\delta))$ times. By the setting of m and τ , the overall query complexity is therefore $O\left(\frac{1}{\varepsilon^2} \log \frac{1}{\varepsilon} \log \frac{1}{\delta}\right)$. \square

5.4 A Non-Adaptive Upper Bound

In this section, we sketch how similar – yet less involved – ideas can be used to derive a non-adaptive upper bound for support size estimation. For simplicity, we describe the algorithm for 2-approximation: adapting it to general $(1 + \varepsilon)$ -approximation is straightforward.

The high-level idea is to perform a simple binary search (instead of the double exponential search from the preceding section) to identify the greatest lower bound on the support size of the form $k = 2^j$. For each guess $k \in \{2, 4, 8 \dots, n\}$, we pick uniformly at random a set $S \subseteq [n]$ of cardinality k , and check whether D_S is uniform using the non-adaptive tester of Chakraborty et al. [CFGM13, Theorem 4.1.2]. If D_S is found to be uniform for all values of k , we return n as our estimate (as the distribution is close to uniform on $[n]$); otherwise, we return n/k , for the smallest k on which D_S was found to be far from uniform. Indeed, D_S can only be far from uniform if S contains points from the support of D , which intuitively only happens if $n/k = \Omega(1)$.

To be more precise, the algorithm proceeds as follows, where $\tau > 0$ is an absolute constant:

```

for all  $k \in \{2, 4, \dots, n\}$  do
  Set a counter  $c_k \leftarrow 0$ .
  for  $m = O(\log \log n)$  times do
    Pick uniformly at random a set  $S \subseteq [n]$  of  $k$  elements.
    Test (non-adaptively) uniformity of  $D_S$  on  $S$ , with the tester of [CFGM13].
    if the tester rejects then increment  $c_k$ .
  end if
end for
if  $c_k > \tau \cdot m$  then return  $\tilde{\omega} \leftarrow \frac{n}{k}$ .
end if
end for
return  $\tilde{\omega} \leftarrow n$ .

```

The query complexity is easily seen to be $\text{poly} \log n$, from the $\tilde{O}(\log n)$ calls to the $\text{poly}(\log n)$ tester of [CFGM13, Theorem 4.1.2]. As for correctness, it follows from the fact that for any set S with mass $D(S) > 0$ which contains at least an η fraction of points outside the support, it holds that D_S is η -far from \mathcal{U}_S .

Acknowledgments. Clément Canonne would like to thank Dana Ron and Rocco Servedio for the many helpful discussions and remarks that influenced the lower bound construction of Section 3.

References

- [ADJ⁺11] Jayadev Acharya, Hirakendu Das, Ashkan Jafarpour, Alon Orlitsky, and Shengjun Pan. Competitive closeness testing. In Sham M. Kakade and Ulrike von Luxburg, editors, *COLT*, volume 19 of *JMLR Proceedings*, pages 47–68. JMLR.org, 2011. 1
- [ADJ⁺12] Jayadev Acharya, Hirakendu Das, Ashkan Jafarpour, Alon Orlitsky, Shengjun Pan, and Ananda Theertha Suresh. Competitive classification and closeness testing. In *Proceedings of 25th COLT*, pages 22.1–22.18, 2012. 1
- [BFF⁺01] Tuğkan Batu, Eldar Fischer, Lance Fortnow, Ravi Kumar, Ronitt Rubinfeld, and Patrick White. Testing random variables for independence and identity. In *Proceedings of FOCS*, pages 442–451, 2001. 1
- [BFR⁺00] Tuğkan Batu, Lance Fortnow, Ronitt Rubinfeld, Warren D. Smith, and Patrick White. Testing that distributions are close. In *Proceedings of FOCS*, pages 189–197, 2000. 1, 5.3
- [BFR⁺10] Tuğkan Batu, Lance Fortnow, Ronitt Rubinfeld, Warren D. Smith, and Patrick White. Testing closeness of discrete distributions. (abs/1009.5397), 2010. This is a long version of [BFR⁺00]. 1.1, 1.1, 1.2
- [BFRV11] Arnab Bhattacharyya, Eldar Fischer, Ronitt Rubinfeld, and Paul Valiant. Testing monotonicity of distributions over general partial orders. In *Proceedings of ITCs*, pages 239–252, 2011. 1

- [BKR04] Tuğkan Batu, Ravi Kumar, and Ronitt Rubinfeld. Sublinear algorithms for testing monotone and unimodal distributions. In *Proceedings of STOC*, pages 381–390, New York, NY, USA, 2004. ACM. [1](#)
- [Can15] Clément L. Canonne. A Survey on Distribution Testing: your data is Big, but is it Blue? *Electronic Colloquium on Computational Complexity (ECCC)*, (TR15-063), April 2015. [1](#)
- [CDVV14] Siu-On Chan, Ilias Diakonikolas, Gregory Valiant, and Paul Valiant. Optimal algorithms for testing closeness of discrete distributions. In *Proceedings of SODA*, pages 1193–1203. Society for Industrial and Applied Mathematics (SIAM), 2014. [1](#), [1.1](#), [1.2](#)
- [CFGM13] Sourav Chakraborty, Eldar Fischer, Yonatan Goldhirsh, and Arie Matsliah. On the power of conditional samples in distribution testing. In *Proceedings of ITCS*, pages 561–580, New York, NY, USA, 2013. ACM. ([document](#)), [1](#), [1.1](#), [1.1](#), [1.2](#), [1.2](#), [1.2](#), [1.2](#), [1.3](#), [1.3](#), [2.1](#), [2.2](#), [3.2.1](#), [5.4](#), [5](#)
- [CR14] Clément L. Canonne and Ronitt Rubinfeld. Testing probability distributions underlying aggregated data. In *Proceedings of ICALP*, pages 283–295, 2014. [1](#)
- [CRS15] Clément L. Canonne, Dana Ron, and Rocco A. Servedio. Testing probability distributions using conditional samples. *SIAM Journal on Computing*, 2015. To appear. Also available on arXiv at [abs/1211.2664](#). ([document](#)), [1](#), [1.1](#), [1.1](#), [2](#), [1.3](#), [2.1](#), [3.1](#)
- [DR96] Devdatt Dubhashi and Desh Ranjan. Balls and Bins: A Study in Negative Dependence. *Random Structures and Algorithms*, 13:99–124, 1996. [3.2.2](#), [4](#)
- [FJO⁺15] Moein Falahatgar, Ashkan Jafarpour, Alon Orlitsky, Venkatadheeraj Pichapathi, and Ananda Theertha Suresh. Faster algorithms for testing under conditional sampling. ([abs/1504.04103](#)), April 2015. ([document](#)), [1.1](#), [1.1](#), [1.2](#), [1.2](#)
- [GMV06] Sudipto Guha, Andrew McGregor, and Suresh Venkatasubramanian. Streaming and sublinear approximation of entropy and information distances. In *Proceedings of SODA*, pages 733–742, Philadelphia, PA, USA, 2006. Society for Industrial and Applied Mathematics (SIAM). [1](#)
- [GR00] Oded Goldreich and Dana Ron. On testing expansion in bounded-degree graphs. Technical Report TR00-020, Electronic Colloquium on Computational Complexity (ECCC), 2000. [1](#), [1.1](#), [1.2](#)
- [ILR12] Piotr Indyk, Reut Levi, and Ronitt Rubinfeld. Approximating and Testing k -Histogram Distributions in Sub-linear Time. In *Proceedings of PODS*, pages 15–22, 2012. [1](#)
- [LRR13] Reut Levi, Dana Ron, and Ronitt Rubinfeld. Testing properties of collections of distributions. *Theory of Computing*, 9(8):295–347, 2013. [1](#)
- [MR95] Rajeev Motwani and Prabhakar Raghavan. *Randomized Algorithms*. Cambridge University Press, New York, NY, USA, 1995. [2.3](#)

- [Pan08] Liam Paninski. A coincidence-based test for uniformity given very sparsely sampled discrete data. *IEEE Transactions on Information Theory*, 54(10):4750–4755, 2008. [1](#), [1.1](#), [1.2](#)
- [RRSS09] Sofya Raskhodnikova, Dana Ron, Amir Shpilka, and Adam Smith. Strong lower bounds for approximating distributions support size and the distinct elements problem. *SIAM Journal on Computing*, 39(3):813–842, 2009. [1.1](#)
- [RS09] Ronitt Rubinfeld and Rocco A. Servedio. Testing monotone high-dimensional distributions. *Random Structures and Algorithms*, 34(1):24–44, January 2009. [1](#)
- [RT14] Dana Ron and Gilad Tsur. The power of an example: Hidden set size approximation using group queries and conditional sampling. *CoRR*, abs/1404.5568, 2014. ([document](#)), [1.2.1](#), [1.3](#)
- [Rub12] Ronitt Rubinfeld. Taming Big Probability Distributions. *XRDS*, 19(1):24–28, September 2012. [1](#)
- [Sto85] L. Stockmeyer. On approximation algorithms for $\#P$. *SIAM Journal on Computing*, 14(4):849–861, 1985. [1.3](#)
- [Sub] List of Open Problems in Sublinear Algorithms: Problem 66. <http://sublinear.info/66>. Bertinoro Workshop on Sublinear Algorithms 2014 (suggested by Eldar Fischer). ([document](#)), [1.1](#), [1.2](#)
- [Val11] Paul Valiant. Testing symmetric properties of distributions. *SIAM Journal on Computing*, 40(6):1927–1968, 2011. [1.1](#)
- [VV10a] Gregory Valiant and Paul Valiant. A CLT and tight lower bounds for estimating entropy. *Electronic Colloquium on Computational Complexity (ECCC)*, (TR10-179), 2010. [1.1](#), [1.2](#), [5.3](#)
- [VV10b] Gregory Valiant and Paul Valiant. Estimating the unseen: A sublinear-sample canonical estimator of distributions. *Electronic Colloquium on Computational Complexity (ECCC)*, (TR10-180), 2010. [5.3](#)
- [VV11] Gregory Valiant and Paul Valiant. The power of linear estimators. In *Proceedings of FOCS*, pages 403–412, October 2011. See also [\[VV10a\]](#) and [\[VV10b\]](#). [1.1](#)
- [VV14] Gregory Valiant and Paul Valiant. An automatic inequality prover and instance optimal identity testing. In *Proceedings of FOCS*, 2014. [1](#), [1.1](#)

A Proof of Lemma 4.2

We now give the proof of the ‘‘Hitting Lemma’’ (Lemma 4.2). We take the logarithms to indicate that a set a_i contributes to C_t with a set of size s if and only if $\log s \in [\log(n/a) - t, \log(n/a) + t]$. Indeed, we can restate the wishful lemma in an additive form. Let $\mathcal{A} = \{a_1, \dots, a_q\}$ be any set of points in $[0, \log n]$. Here the new a_i ’s are related to the original a_i ’s, i.e. they are $\log(n/a_i)$. For a point $x \in [0, n]$, let ℓ_j^x (r_j^x respectively) denote the distance of x from the j th point to its left (right respectively) from the set \mathcal{A} . (So that x corresponds to $\log s$.) Let

$$q_x \stackrel{\text{def}}{=} \min_j \min \left\{ \frac{\ell_j^x}{j}, \frac{r_j^x}{j} \right\}.$$

For a constant $c > 0$, let S_c be the set of all points x such that $q_x < c$. A point x violates the hitting lemma if for some j ,

$$\min \left\{ \frac{\ell_j^x}{j}, \frac{r_j^x}{j} \right\} < \log \beta \cdot \frac{2}{100}.$$

Therefore, we need to bound the size of S_c for $c = 2 \log \beta / 100$. We do this by proving the following. (This will imply the Hitting Lemma, as q is at most $\log n / 200 \log \beta$ and the set of possible values for x has size $\frac{1}{2} \log n$.)

Lemma A.1. *Let $|S_c|$ be the measure of S_c . Then $|S_c| \leq 2cq$.*

Proof of Lemma A.1. We consider the set of points in $S_{c,\ell} \subset S_c$ that satisfy $\ell_j^x/j < c$ for some j , and show that their measure is at most cq . An identical bound holds for the set of points of S_c for which $r_j^x/j < c$. Let $S_{c,\ell}^i \subset S_{c,\ell}$ be the set of points in $S_{c,\ell}$ that satisfy $m_x < c$ with respect to the set a_1, \dots, a_i . We will show by induction that $|S_{c,\ell}^i| < ci$.

For the first point a_1 , the set $S_{c,\ell}^1 = [a_1, a_1 + c]$. Suppose by induction that $|S_{c,\ell}^i| < ci$. Let x_i be the right-most point in the set $S_{c,\ell}^i$. Then it is clear that $x_i > a_i$, in fact $x_i \geq a_i + c$. Furthermore, either $x_i = n$, or $\ell_j^{x_i}/j = c$ for some j . Moreover, we claim that $[a_i, x_i] \in S_{c,\ell}^i$. Indeed, for the same j that $\ell_j^{x_i}/j < c$, all points in $[a_i, x_i]$ satisfy the condition. If $x_i = n$, then the result holds trivially. We therefore consider the point a_{i+1} and prove the inductive step for $x_i < n$. There are two cases:

If $a_{i+1} \geq x_i$: In this case, $S_{c,\ell}^{i+1} = S_{c,\ell}^i \cup [a_{i+1}, x_{i+1}]$. We have to show that $x_{i+1} \leq a_{i+1} + c$. Suppose to the contrary that $x_{i+1} > a_{i+1} + c \geq x_i + c$. Then there is a point a_h for $h \leq i$, such that $\frac{x_{i+1} - a_h}{i+2-h} < c$, and then $\frac{a_{i+1} + c - a_h}{i+2-h} < c$, so that

$$\frac{a_{i+1} - a_h}{i+1-h} < c,$$

however, this implies that $a_{i+1} \in S_{c,\ell}^i$, contradicting the assumption of this case.

If $a_{i+1} < x_i$: In this case, $S_{c,\ell}^{i+1} = S_{c,\ell}^i \cup [x_i, x_{i+1}]$. We have to show that $x_{i+1} \leq x_i + c$. Suppose on the contrary that $x_{i+1} > x_i + c > a_{i+1} + c$. Suppose h be the index such that $\frac{x_{i+1} - a_h}{i+2-h} < c$, and therefore, $\frac{x_i + c - a_h}{i+2-h} < c$, implying that

$$\frac{x_i - a_h}{i+1-h} < c,$$

contradicting that x_i is the rightmost point of $S_{c,\ell}^i$. □