

# Lower bounds for testing function isomorphism

Eric Blais  
Computer Science Department  
Carnegie Mellon University  
Pittsburgh, PA 15213  
eblais@cs.cmu.edu

Ryan O’Donnell\*  
Computer Science Department  
Carnegie Mellon University  
Pittsburgh, PA 15213  
odonnell@cs.cmu.edu

**Abstract**—We prove new lower bounds in the area of property testing of boolean functions. Specifically, we study the problem of testing whether a boolean function  $f$  is isomorphic to a fixed function  $g$  (i.e., is equal to  $g$  up to permutation of the input variables). The analogous problem for testing graphs was solved by Fischer in 2005. The setting of boolean functions, however, appears to be more difficult, and no progress has been made since the initial study of the problem by Fischer et al. in 2004.

Our first result shows that any non-adaptive algorithm for testing isomorphism to a function that “strongly” depends on  $k$  variables requires  $\log k - O(1)$  queries (assuming  $k/n$  is bounded away from 1). This lower bound affirms and strengthens a conjecture appearing in the 2004 work of Fischer et al. Its proof relies on total variation bounds between hypergeometric distributions which may be of independent interest.

Our second result concerns the simplest interesting case not covered by our first result: non-adaptively testing isomorphism to the Majority function on  $k$  variables. Here we show that  $\Omega(k^{1/12})$  queries are necessary (again assuming  $k/n$  is bounded away from 1). The proof of this result relies on recently developed multidimensional invariance principle tools.

**Keywords**—Boolean functions, property testing, lower bounds

## I. INTRODUCTION

This paper is concerned with the field of property testing for boolean functions. Let us recall the standard framework, as originally introduced by Rubinfeld and Sudan [17].

**Definition 1.1.** Let  $\mathcal{P}$  be a class of functions  $\{0, 1\}^n \rightarrow \{0, 1\}$ . We say that a randomized query algorithm  $\mathcal{T}$  with black-box access to an unknown function  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  is an  $(\epsilon, q)$ -tester for  $\mathcal{P}$  if it makes at most  $q$  queries to  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  and then:

- Accepts with probability at least  $2/3$  when  $f$  is in  $\mathcal{P}$ ; and,
- Rejects with probability at least  $2/3$  when  $f$  is  $\epsilon$ -far from every function  $f' \in \mathcal{P}$ .

\* Supported by NSF grants CCF-0747250 and CCF-0915893, BSF grant 2008477, and Sloan and Okawa fellowships.

Here we say that  $f$  and  $f'$  are  $\epsilon$ -far if they differ on at least an  $\epsilon$  fraction of the inputs in  $\{0, 1\}^n$ , and are  $\epsilon$ -close otherwise. When the algorithm chooses all of its queries in advance it is *non-adaptive*; otherwise we say it is *adaptive*.

**Definition 1.2.** For a fixed  $\epsilon > 0$  and choice of adaptivity, the *query complexity* of  $\mathcal{P}$  is the minimum value of  $q$  for which there is an  $(\epsilon, q)$ -tester. Following standard conventions, when the query complexity  $q$  is independent of  $n$  for every  $\epsilon > 0$ , we say that  $\mathcal{P}$  is *easy to test*; otherwise we say that it is *hard to test*.<sup>1</sup> This notion is independent of the choice of adaptivity, since non-adaptive query complexity can be (exponentially) bounded in terms of adaptive query complexity.

The last five years have seen great strides in understanding the testability of *graph properties*. This is the special case in which  $f : \{0, 1\}^{\binom{V}{2}} \rightarrow \{0, 1\}$  encodes the adjacency matrix of a graph on vertex set  $V$ , and  $\mathcal{P}$  is a property that is closed under graph symmetries; i.e., permutations of  $V$ . Indeed, the works [2], [3], [1], [4] have to a large extent characterized the testability of graph properties.

However the characterization problem for general boolean functions is very far from understood and remains a long-standing open problem. In this paper we revisit a major subproblem, introduced in the early work of Fischer, Kindler, Ron, Safra, and Samorodnitsky [7]: the difficulty of testing isomorphism to a function given in advance.

### A. Testing $g$ -isomorphism

Two boolean functions  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  and  $g : \{0, 1\}^n \rightarrow \{0, 1\}$  are said to be *isomorphic* to each other if they are identical up to reordering input variables. More precisely, we say that  $f$  and  $g$  are isomorphic to each other if there is a permutation  $\sigma$  on  $[n]$  such that for every  $x = (x_1, x_2, \dots, x_n) \in \{0, 1\}^n$ ,  $f(x_1, x_2, \dots, x_n) = g(x_{\sigma(1)}, x_{\sigma(2)}, \dots, x_{\sigma(n)})$ .

For each function  $g : \{0, 1\}^n \rightarrow \{0, 1\}$ , we let  $\mathcal{P}_g$  denote the class of all functions isomorphic to  $g$ . This gives a natural

<sup>1</sup>Formally speaking, this makes sense only for a family of properties  $(\mathcal{P}_n)_n$ , one for each input length.

testing problem: testing whether an unknown function  $f$  is isomorphic to the known function  $g$ . This is called the problem of *testing  $g$ -isomorphism*. Fischer et al. [7] proposed the following question, a significant component of the research program to characterize testability of boolean functions:

**Research goal:** Classify all boolean functions  $g$  according to whether testing  $g$ -isomorphism is easy or hard.

We remark that Fischer [6] solved the analogous problem for *graph properties* soon after; see also [8]. But he called the general case of boolean functions “rather hard”, and indeed the authors are not aware of any additional progress specifically on this problem.

In this paper we make progress towards a characterization by showing hardness of testing  $g$ -isomorphism for a large class of functions  $g$ .

### B. Prior work

When the problem of testing function isomorphism was first raised by Fischer et al. [7], a few simple cases were already well understood. First, it is easy to see that for any totally symmetric function  $g : \{0, 1\}^n \rightarrow \{0, 1\}$ , testing  $g$ -isomorphism is easy – no other functions are isomorphic to  $g$ , and testing function *identity* requires only  $O(1/\epsilon)$  queries.

Another instance of the  $g$ -isomorphism testing problem that is well understood is one where  $g(x) = x_i$  for some  $i \in \{1, \dots, n\}$ . Then the  $g$ -isomorphism problem is equivalent to the well-studied *dictatorship* testing problem at the heart of PCP constructions (first studied in [5]). The query complexity of the dictatorship testing problem is  $O(1/\epsilon)$ , so this special case of the function isomorphism problem is also easy. Parnas, Ron, and Samorodnitsky [15] also showed that testing  $g$ -isomorphism is easy when  $g$  is an AND function on any number of variables.

The paper of Fischer et al. [7] introduced a strong new upper bound: they showed that for any  $k$ , if  $g$  is a  $k$ -junta – meaning that  $g$  depends on at most  $k$  variables – then it is possible to  $\epsilon$ -test  $g$ -isomorphism (non-adaptively, even) with  $\text{poly}(k/\epsilon)$  queries. Therefore, testing isomorphism to any  $O(1)$ -junta is easy.

Regarding hardness results, prior to this work the only known lower bound for testing  $g$ -isomorphism was for the case of  $g$  being the  $\text{Parity}_k$  function, where  $\text{Parity}_k : \{0, 1\}^n \rightarrow \{0, 1\}$  is defined by  $\text{Parity}_k(x) = x_1 \oplus x_2 \oplus \dots \oplus x_k$ . Fischer et al. [7] showed that when  $k \leq o(\sqrt{n})$ ,  $\epsilon$ -testing  $\text{Parity}_k$ -isomorphism non-adaptively requires  $\tilde{\Omega}(\sqrt{k}/\epsilon)$  queries. This result implies that testing  $\text{Parity}_k$ -isomorphism is hard for  $\omega(1) \leq k \leq o(\sqrt{n})$ .

### C. Our results

It seems clear that more work needs to be done on the hardness side of testing  $g$ -isomorphism. A first direction

would be to investigate the following conjecture, stated in [7]: “If  $n$  is sufficiently large compared to  $k$ , and  $g$  is a  $k$ -junta which is  $\epsilon$ -far from all  $(k-1)$ -juntas, then  $\epsilon$ -testing  $g$ -isomorphism requires  $\omega_k(1)$  queries.”

Our first main result affirms and significantly strengthens this conjecture:

**Theorem 1.3.** *Let  $g : \{0, 1\}^n \rightarrow \{0, 1\}$  be a  $k$ -junta which is  $\epsilon$ -far from being a  $(k-e)$ -junta for some  $e \geq 1$ . Then any non-adaptive  $\epsilon$ -tester for  $g$ -isomorphism must make at least  $\log_2(k'/e^2) - O(1)$  queries, where  $k' = \min(k, n-k)$ .*

Qualitatively, Theorem 1.3 implies that testing  $g$ -isomorphism for  $k$ -juntas is hard whenever  $\omega(1) \leq k \leq n - \omega(1)$ , provided the  $k$ -junta  $g$  is far from all juntas on  $k - o(\sqrt{k})$  variables. We discuss the possibility of improving this theorem in Section V.

As one very special case, Theorem 1.3 extends the Fischer et al. hardness result for testing  $\text{Parity}_k$ -isomorphism to all  $\omega(1) \leq k \leq n - \omega(1)$  (albeit with a worse query complexity lower bound). It is important to note that the restriction  $k \leq n - \omega(1)$  is not an artifact of our proof, but rather is inherent. This is because, e.g., testing  $\text{Parity}_k$ -isomorphism when  $k = n - O(1)$  is *easy*: if the tester XORs each query response with the parity of all bits in the query, it reduces testing  $\text{Parity}_k$ -isomorphism to testing  $\text{Parity}_{n-k}$ -isomorphism. Thus testing  $\text{Parity}_{n-O(1)}$ -isomorphism is easy, by the junta isomorphism-testing result of Fischer et al.

Our Theorem 1.3 is only useful for  $k$ -juntas  $g$  that are far from being  $(k - o(\sqrt{k}))$ -juntas. Perhaps the most natural case not covered by this theorem is that of the majority function on  $k$  variables,  $\text{Maj}_k$ . A straightforward calculation shows that for all  $\delta > 0$ , the function  $\text{Maj}_k$  is  $o_\delta(1)$ -close to being a junta on  $k - \delta k$  variables: namely  $\text{Maj}_{k-\delta k}$ . Nevertheless, the second main result in our paper proves a strong hardness result for testing  $\text{Maj}_k$ -isomorphism:

**Theorem 1.4.** *For every constant  $\delta > 0$ , there exists a constant  $\epsilon > 0$  such that the following holds: Assuming  $1/\epsilon \leq k \leq (1 - \delta)n$ , any non-adaptive algorithm for  $\epsilon$ -testing  $\text{Maj}_k$ -isomorphism must make at least  $\Omega((\delta k)^{1/12})$  queries.*

Qualitatively, Theorem 1.4 implies that testing  $\text{Maj}_k$ -isomorphism is hard for every  $\omega(1) \leq k \leq n - \Omega(n)$ . Again, this range is optimal: the upper bound cannot be improved because as mentioned,  $\text{Maj}_{n-o(n)}$  is  $o(1)$ -close to  $\text{Maj}_n$ ; thus we can test  $\text{Maj}_{n-o(n)}$ -isomorphism using the easy  $\text{Maj}_n$ -isomorphism tester that follows from  $\text{Maj}_n$  being totally symmetric.

In this extended abstract, we prove Theorem 1.4 only in the case where  $\delta = 1/4$  (and assuming  $k$  is divisible by 3). The few tedious technical modifications needed to handle smaller values of  $\delta$  are deferred to the full version of the

article. I.e., we prove:

**Theorem 1.5.** *There is a universal  $\epsilon_0 > 0$  such that whenever  $k \leq (3/4)n$  (and is divisible by 3), any non-adaptive algorithm for  $\epsilon_0$ -testing  $\text{Maj}_k$ -isomorphism must make at least  $\Omega(k^{1/12})$  queries.*

We remark that the bound in Theorem 1.5 is exponentially stronger than the one in Theorem 1.3: it shows that any non-adaptive algorithm for testing  $\text{Maj}_k$ -isomorphism must make a number of queries *polynomial* in  $k$ . This bound is optimal (up to the right value of the exponent), since Fischer et al.’s [7] upper bound on the query complexity of testing isomorphism to  $k$ -juntas implies that testing isomorphism to the  $\text{Maj}_k$  function can be done with  $\text{poly}(k/\epsilon)$  queries.

The result of Theorem 1.4 is also interesting in light of the recent results of Matulef et al. on testing halfspaces: they showed that testing the class of halfspaces is easy [11], but testing a natural subclass of halfspaces – the class of  $\pm 1$ -weight halfspaces – is hard [12]. More precisely, they showed a non-adaptive query lower bound of  $\Omega(\log n)$  for this class. Our Theorem 1.4 gives an exponentially improved lower bound for a similar subclass of halfspaces:  $\Omega(n^{1/12})$  queries for the class of majority functions on, say,  $n/2$  variables.

In light of Fischer [6]’s solution to the isomorphism testing problem for graph properties – i.e., boolean functions with a certain high degree of symmetry – we believe that characterizing testability of  $g$ -isomorphism for *symmetric*  $k$ -juntas is an approachable first step. Theorem 1.5 represents progress in this direction. We remark that our method of proving Theorem 1.5 can be extended to handle certain other symmetric  $k$ -juntas. However the general case of symmetric  $k$ -juntas  $g$  has some unexpected tricky aspects to it, which we discuss in Section V.

#### D. Our techniques

The proofs of Theorems 1.3 and 1.5 both use the standard approach for proving lower bounds in property testing: Yao’s Minimax Principle [18]. That is, we prove both theorems by introducing distributions  $\mathcal{F}_{\text{yes}}$  and  $\mathcal{F}_{\text{no}}$  on functions that should be accepted and rejected, respectively, by algorithms testing  $g$ -isomorphism, and show a lower bound on the number of queries required by any *deterministic* testing algorithm. The main technical contribution of this research is in the design and the analysis of the distributions  $\mathcal{F}_{\text{yes}}$  and  $\mathcal{F}_{\text{no}}$ .

The main challenge in proving Theorem 1.3 is that the lower bound applies to a very general class of functions  $g$ . To prove the theorem we need to design distributions that work *without using any structural properties of the function  $g$  being tested*. The key to doing this involves analyzing the statistical (total variation) distance between two *multivariate*

*hypergeometric distributions*. What follows is the main lemma we need; it may be of independent interest:

**Lemma 1.6.** *Suppose  $X \sim \text{Hyp}(n, r, k)$  and  $Y \sim \text{Hyp}(n, r, \ell + e)$ , where  $\text{Hyp}(n, r, \ell)$  denotes the (univariate) hypergeometric distribution: i.e., the number of red balls drawn when selecting  $\ell$  balls randomly without replacement from an urn containing  $n$  balls,  $r$  of which are red. Then*

$$d_{\text{TV}}(X, Y) \leq .01$$

*provided  $(1 - \frac{r}{n}) \min(\ell, n - \ell) \geq Ce^2$ , where  $C$  is a universal constant. Here  $d_{\text{TV}}(\cdot, \cdot)$  denotes total variation distance.*

In Section III-E we comment on why the somewhat complicated hypothesis on  $r$ ,  $n$ ,  $\ell$ , and  $e$  is necessary. The proof of Theorem 1.3, as well as a more complete discussion of the techniques it requires, is presented in Section III.

The challenge in proving Theorem 1.5 is fundamentally different. For Theorem 1.3, we know that the  $k$ -junta  $g$  is far from all  $(k - 1)$ -juntas, say, which means it is okay for our  $\mathcal{F}_{\text{no}}$  functions to be “small tweaks” to  $g$ . However when  $g = \text{Maj}_k$ , small tweaks result in functions that are still close to  $\text{Maj}_k$ . Thus our  $\mathcal{F}_{\text{no}}$  functions must be somewhat drastically changed from  $\text{Maj}_k$ , yet still “look like”  $\text{Maj}_k$ . We arrange for this by making the  $\mathcal{F}_{\text{no}}$  functions very carefully constructed *weighted* majority functions on  $\frac{4}{3}k$  coordinates. To show that such functions still “look like”  $\text{Maj}_k$ , we use recently developed multidimensional invariance principles [13], [9]. However these need to be adapted to the case of sums of random vectors which are *not independent*, but rather are drawn without replacement from a fixed pool of random vectors.

## II. PRELIMINARIES AND DEFINITIONS

### A. Probability theory

Given two random variables  $X, Y$  defined on a common discrete sample space  $\Omega$ , the *total variation* distance between  $X$  and  $Y$  is

$$d_{\text{TV}}(X, Y) = \frac{1}{2} \sum_{\omega \in \Omega} |\Pr[X = \omega] - \Pr[Y = \omega]|.$$

When  $X = (X_1, \dots, X_n) \in \mathbb{R}^n$  is a random vector and we let  $\mu_i = \mathbf{E}[X_i]$  for  $i = 1, \dots, n$ , then the *mean* of  $X$  is the vector  $\mathbf{E}[X] = (\mu_1, \dots, \mu_n)$  and the *covariance matrix* of  $X$  is the  $n \times n$  matrix  $\mathbf{Cov}[X]$  whose  $(i, j)$ th entry is defined by

$$\mathbf{Cov}[X]_{i,j} = \mathbf{E}[(X_i - \mu_i)(X_j - \mu_j)] = \mathbf{E}[X_i X_j] - \mu_i \mu_j.$$

The function  $\Phi : \mathbb{R} \rightarrow [0, 1]$  is used throughout this paper to denote the cumulative density function of the standard normal distribution  $\mathcal{N}(0, 1)$ .

We write  $\text{Hyp}(n, m, k)$  to represent the hypergeometric distribution – the distribution on the number  $t$  of red balls

drawn when  $k$  balls are drawn without replacement from a set of  $n$  balls,  $m$  of which are red.

The following theorem of Höglund [10] provides useful bounds on the normal approximation of hypergeometric distributions.

**Höglund Theorem.** *Let  $S = X_1 + \dots + X_n$ , where  $X_1, \dots, X_n$  are chosen uniformly at random without replacement from  $A = \{x_1, \dots, x_N\}$ . Then for all  $t \in \mathbb{R}$ ,*

$$\left| \Pr[S \leq t] - \Phi\left(\frac{t - n\mu}{\sigma\sqrt{n(1-n/N)}}\right) \right| \leq C \cdot \frac{\sum_{i=1}^N |x_i - \mu|^3 / N}{\sigma^3 \sqrt{n(1-n/N)}},$$

where  $\mu = \sum_{i=1}^N \frac{x_i}{N}$ ,  $\sigma^2 = \sum_{i=1}^N \frac{(x_i - \mu)^2}{N}$ , and  $C$  is an absolute constant.

Höglund's Theorem yields the following anti-concentration property of hypergeometric distributions.

**Corollary 2.1.** *Let  $S \sim \text{Hyp}(n, m, k)$  and let  $\sigma^2 = \frac{km}{n} \left(1 - \frac{m}{n}\right) \left(1 - \frac{k}{n}\right)$ . Then for every  $t \geq 0$ ,*

$$\Pr[S = t] \leq \frac{C}{\sigma}$$

where  $C$  is an absolute constant.

*Proof:* Let  $A$  be a set containing  $m$  ones and  $n - m$  zeros. The random variable  $S = X_1 + \dots + X_k$  obtained by choosing  $k$  variables uniformly at random without replacement from  $A$  follows the  $\text{Hyp}(n, m, k)$  distributions. Fixing  $\mu = km/n$ , Höglund's Theorem implies that for any  $t \in \mathbb{R}$ ,

$$\left| \Pr[S \leq t] - \Phi\left(\frac{t - \mu}{\sigma}\right) \right| \leq c_0 \cdot \frac{\sum |x_i - m/n|^3}{\sum (x_i - m/n)^2} \cdot \frac{1}{\sigma} \leq \frac{c_0}{\sigma}.$$

Since  $\Pr[S = t] = \Pr[S \leq t] - \Pr[S \leq t - 1]$ , we have

$$\Pr[S = t] \leq \Phi\left(\frac{t - \mu}{\sigma}\right) - \Phi\left(\frac{t - 1 - \mu}{\sigma}\right) + 2\frac{c_0}{\sigma} \leq \frac{C}{\sigma},$$

where  $C = 2c_0 + \sqrt{2/\pi}$ .  $\blacksquare$

### B. Influence

The *influence* of the  $i$ th variable in the function  $f : \{0, 1\}^n \rightarrow \{0, 1\}$  is  $\text{Inf}_f(i) = \Pr_x[f(x) \neq f(x^{(i)})]$ , where the probability is taken over the uniform distribution of  $x \in \{0, 1\}^n$  and  $x^{(i)}$  is the input formed by flipping the value of the  $i$ th bit in  $x$ .

When  $\text{Inf}_f(i) > 0$ , we say that the  $i$ th variable is *relevant* in  $f$ . A function that contains at most  $k$  relevant variables is called a  $k$ -*junta*.

## III. GENERAL LOWER BOUND

In this section, we prove Theorem 1.3.

**Theorem 1.3 (Restated).** *Let  $g : \{0, 1\}^n \rightarrow \{0, 1\}$  be a  $k$ -junta which is  $\epsilon$ -far from being a  $(k - e)$ -junta for some  $e \geq 1$ . Then any non-adaptive  $\epsilon$ -tester for  $g$ -isomorphism must make at least  $\log_2(k'/e^2) - O(1)$  queries, where  $k' = \min(k, n - k)$ .*

On first reading, the reader is encouraged to focus on the simplest case, where  $e = 1$ . In this case, Theorem 1.3 affirms a conjecture stated in [7], as it shows that for any  $\omega(1) \leq k \leq n - \omega(1)$ , when  $g$  is a  $k$ -junta and is  $\epsilon$ -far from being a  $(k - 1)$ -junta for some  $\epsilon > 0$ , then testing  $g$ -isomorphism is hard.

We now begin the proof of Theorem 1.3. Let  $g, n, k, k', e$ , and  $\epsilon$  be as in the statement of the theorem. Without loss of generality, we may assume that the  $k$  relevant coordinates for  $g$  are  $[k] = \{1, 2, \dots, k\}$ . We write  $g_{\text{core}} : \{0, 1\}^k \rightarrow \{0, 1\}$  for the restriction of  $g$  to these coordinates.

As is standard in property testing lower bounds, the proof of Theorem 1.3 uses Yao's Minimax Principle [18]. Specifically, we construct two probability distributions  $\mathcal{F}_{\text{yes}}$  and  $\mathcal{F}_{\text{no}}$  over functions isomorphic to  $g$  and functions  $\epsilon$ -far from being isomorphic to  $g$ , respectively. We then show that any deterministic non-adaptive algorithm making  $\ll \log_2(k'/e^2)$  queries cannot distinguish with probability at least  $1/3$  between functions drawn from  $\mathcal{F}_{\text{yes}}$  or from  $\mathcal{F}_{\text{no}}$ .

### A. The distributions $\mathcal{F}_{\text{yes}}$ and $\mathcal{F}_{\text{no}}$

We define  $\mathcal{F}_{\text{yes}}$  in the most natural way, by randomly embedding  $g_{\text{core}}$  into  $[n]$ . More precisely, to obtain a draw  $f_{\text{yes}} \sim \mathcal{F}_{\text{yes}}$ , we first choose a uniformly random subset  $J \subseteq [n]$  of cardinality  $k$ . Next, we choose a uniformly random bijection  $\sigma : [k] \rightarrow J$ . Finally, we define  $f_{\text{yes}}(x) = g_{\text{core}}(x_{\sigma(1)}, \dots, x_{\sigma(k)})$ . It is clear that every such  $f_{\text{yes}}$  is isomorphic to  $g$ .

As for  $\mathcal{F}_{\text{no}}$ , we define a draw  $f_{\text{no}} \sim \mathcal{F}_{\text{no}}$  as follows: First, we choose a uniformly random subset  $J \subseteq [n]$  of cardinality  $k - e$ . Next, we choose a uniformly random map  $\sigma : [k] \rightarrow J$  from among those satisfying the following property: there is one  $j_1 \in J$  with  $e + 1$  preimages under  $\sigma$ , and the remaining  $j \in J \setminus \{j_1\}$  have a unique preimage. Finally, we define  $f_{\text{no}}(x) = g_{\text{core}}(x_{\sigma(1)}, \dots, x_{\sigma(k)})$ . Each such  $f_{\text{no}}$  only depends on the coordinates  $J$  and hence is a  $(k - e)$ -junta. Thus by the assumption in Theorem 1.3, each such  $f_{\text{no}}$  is indeed  $\epsilon$ -far from being isomorphic to  $g$ .

To prove Theorem 1.3, it suffices to prove the following:

**Theorem 3.1.** *Let  $\mathcal{T}$  be any deterministic non-adaptive  $q$ -query testing algorithm for functions  $\{0, 1\}^n \rightarrow \{0, 1\}$ . Then*

$$\left| \Pr_{f_{\text{yes}} \sim \mathcal{F}_{\text{yes}}} [\mathcal{T} \text{ acc. } f_{\text{yes}}] - \Pr_{f_{\text{no}} \sim \mathcal{F}_{\text{no}}} [\mathcal{T} \text{ acc. } f_{\text{no}}] \right| \leq O\left(\frac{2^q e^2}{k'}\right) + .01.$$

Note that if  $q < \log_2(k'/e^2) - c_0$  for a sufficiently large constant  $c_0$ , then the upper bound in this theorem is at most  $1/3$ . From this we deduce Theorem 1.3 immediately using Yao's Minimax Principle.

### B. Distance between multivariate hypergeometrics

The typical way to prove a property testing bound such as Theorem 3.1 is as follows. First, we write the  $q$  queries

of tester  $\mathcal{T}$  as  $x^1, \dots, x^q \in \{0, 1\}^n$ . We then introduce the *response vector* random variables  $R_{\text{yes}}$  and  $R_{\text{no}}$ . Here  $R_{\text{yes}} \in \{0, 1\}^q$  is defined by drawing  $f_{\text{yes}} \sim \mathcal{F}_{\text{yes}}$  and letting  $R_{\text{yes}} = \langle f_{\text{yes}}(x^1), \dots, f_{\text{yes}}(x^q) \rangle$ , and  $R_{\text{no}}$  is defined analogously. Finally, we show that

$$d_{\text{TV}}(R_{\text{yes}}, R_{\text{no}}) \leq 2^q \cdot \frac{O(e^2)}{k'} + .01. \quad (1)$$

We will in fact prove a stronger statement. To understand it, let's reconsider the complete random processes  $\mathcal{P}_{\text{yes}}$  and  $\mathcal{P}_{\text{no}}$  by which the response vectors  $R_{\text{yes}}$  and  $R_{\text{no}}$  are generated. We begin by focusing on the “yes” process,  $\mathcal{P}_{\text{yes}}$ .

Given the tester  $\mathcal{T}$ 's queries  $x^1, \dots, x^q \in \{0, 1\}^n$ , we think of them as row vectors and arrange them into a  $q \times n$  *query matrix*  $Q$ . We will be especially interested in the *columns* of this matrix  $Q$ , the  $j$ th column consisting of the  $j$ th bits of all the query strings. Abstractly, we define the set of all possible column (types)

$$\mathcal{C} = \{0, 1\}^q.$$

Since  $|\mathcal{C}| = 2^q \ll n$ , some columns will occur many times in the matrix  $Q$ . In fact, we will think of the query matrix  $Q$  as being an ordered *multiset* of columns from  $\mathcal{C}$ .

Recalling the definition of  $\mathcal{F}_{\text{yes}}$ , we think of the first step of  $\mathcal{P}_{\text{yes}}$  as choosing  $k$  column indices  $j_1, \dots, j_k$  randomly and without replacement from  $[n]$ . We next extract columns  $j_1, \dots, j_k$  from  $Q$ . We view this as a multiset of columns, and call it the *argument multiset*  $S_{\text{yes}}$ . Next, we randomly order the columns in  $S_{\text{yes}}$ , forming a  $q \times k$  *argument matrix*  $A_{\text{yes}}$ . Finally, we produce the response vector  $R_{\text{yes}}$  by applying  $g_{\text{core}}$  to the argument matrix, row-wise.

The reader can easily verify this process  $\mathcal{P}_{\text{yes}}$  generates the correct distribution on the response vector random variable  $R_{\text{yes}}$ .

The “no” process  $\mathcal{P}_{\text{no}}$  is very similar, differing only in the way it generates the argument multiset from the query matrix. Recalling the definition of  $\mathcal{F}_{\text{no}}$ , we think of  $\mathcal{P}_{\text{no}}$  as forming the argument multiset  $S_{\text{no}}$  by choosing  $\ell = k - e$  random columns from  $Q$  without replacement, and including an *additional*  $e$  copies of the first-chosen column. The process  $\mathcal{P}_{\text{no}}$  then forms the argument matrix  $A_{\text{no}}$  by again randomly ordering the columns in the argument multiset, and finally produces the response vector  $R_{\text{no}}$  again by applying  $g_{\text{core}}$  to  $A_{\text{no}}$ , row-wise. The reader can again easily verify that  $\mathcal{P}_{\text{no}}$  generates the correct distribution on  $R_{\text{no}}$ .

Because the processes are identical after the argument multiset is formed, a coupling argument immediately implies that

$$d_{\text{TV}}(R_{\text{yes}}, R_{\text{no}}) \leq d_{\text{TV}}(S_{\text{yes}}, S_{\text{no}}). \quad (2)$$

This inequality can be extremely lossy, depending on the function  $g_{\text{core}}$ . However, since Theorem 1.3 applies for an extremely broad range of functions, we are almost forced to

design a proof of Theorem 3.1 that *uses no properties of the function*  $g_{\text{core}}$ . That is, in the absence of additional restrictions on the class of functions considered, there is no obvious way to bound  $d_{\text{TV}}(R_{\text{yes}}, R_{\text{no}})$  except by  $d_{\text{TV}}(S_{\text{yes}}, S_{\text{no}})$ .

Letting  $\mathcal{S}_{\text{yes}}$  denote the subprocess of  $\mathcal{P}_{\text{yes}}$  generating  $S_{\text{yes}}$ , and similarly for  $\mathcal{S}_{\text{no}}$ , we have reduced proving (1), and hence Theorem 1.3, to the following:

**Theorem 3.2.** *For  $S_{\text{yes}} \sim \mathcal{S}_{\text{yes}}, S_{\text{no}} \sim \mathcal{S}_{\text{no}}$ , we have*

$$d_{\text{TV}}(S_{\text{yes}}, S_{\text{no}}) \leq |\mathcal{C}| \cdot \frac{O(e^2)}{\min(k, n - k)} + .01.$$

The reader can see now why our query complexity lower bound in Theorem 1.3 is only logarithmic; we have  $|\mathcal{C}| = 2^q$  competing against  $\frac{1}{k}$  in the above bound. Indeed, we can never prove a better-than-logarithmic lower bound if our proof only involves showing statistical closeness of the argument multisets  $S_{\text{yes}}$  and  $S_{\text{no}}$ . To see this, suppose  $k = n/2$ , so  $n - k = n/2$  as well. Then if  $2^q \gg n/2$ , it is possible that every column in the query matrix is unique. In this case, the total variation distance between argument multisets  $S_{\text{yes}}$  and  $S_{\text{no}}$  will be 1 even in the case  $e = 1$ , because  $S_{\text{yes}}$  will always consist of unique columns, whereas  $S_{\text{no}}$  will always have one column duplicated.

Notice that the ordering of the columns in the query matrix  $Q$  has proven to be unimportant; we can think of  $Q$  simply as an unordered multiset of columns from  $\mathcal{C}$ . Thus Theorem 3.2 is really a statement about the total variation distance between certain multivariate hypergeometric random variables. Specifically, for each column  $\mathfrak{c} \in \mathcal{C}$ , let  $m(\mathfrak{c})$  denote the number of copies of  $\mathfrak{c}$  in  $Q$ . In process  $\mathcal{S}_{\text{yes}}$ , we choose  $k$  random columns from  $Q$  without replacement and count the number of copies of each column (type) in the draw. Process  $\mathcal{S}_{\text{no}}$  is similar, except we choose  $\ell$  random columns from  $Q$  without replacement, and count an extra  $e$  copies of the first-drawn column.

### C. Reduction of Theorem 3.2 to two lemmas

This preceding discussion motivates the following notation:

**Definition 3.3.** Given integers  $N, e \geq 1, M, L \geq 0$ , with  $M, L + e \leq N$ , we define  $\lambda_{N, M, L}(e) = d_{\text{TV}}(X, Y)$ , where  $X \sim \text{Hyp}(N, M, L + e)$  and  $Y \sim \text{Hyp}(N, M, L) + e$ .

The proof of Theorem 3.2 relies on the following two lemmas. The first lemma is relatively straightforward, and relates the distance between  $\mathcal{S}_{\text{yes}}$  and  $\mathcal{S}_{\text{no}}$  to the total variation distance between hypergeometric distributions.

**Lemma 3.4.**

$$d_{\text{TV}}(S_{\text{yes}}, S_{\text{no}}) \leq \sum_{\mathfrak{c} \in \mathcal{C}: m(\mathfrak{c}) \neq 0} \frac{m(\mathfrak{c})}{n} \cdot \lambda_{n-1, m(\mathfrak{c})-1, \ell-1}(e).$$

The second lemma is a total variation distance bound between (univariate) hypergeometric random variables which may be of independent interest.

**Lemma 3.5.** *There is a universal constant  $2 \leq \kappa < \infty$  such that for any  $N, M, L$ , if  $L' = \min(L, N - L)$  satisfies  $\frac{ML'}{N} \geq \kappa e^2$ , then  $\lambda_{N,M,L}(e) \leq .01$ .*

This lemma is in fact identical to the key Lemma 1.6: to see this, one only needs to replace  $M$  with  $r = N - M$  and use the obvious fact that  $\text{Hyp}(N, N - M, L)$  is the same distribution as  $L - \text{Hyp}(N, M, L)$ .

We briefly comment on why the hypothesis  $\frac{ML'}{N} \gg e^2$  is necessary to show that  $\text{Hyp}(N, M, L + e)$  and  $\text{Hyp}(N, M, L) + e$  are close in total variation distance. For simplicity, first suppose that  $e = 1$ . It is necessary that  $\frac{ML'}{N} \gg 1$ ; this quantity is the mean of  $\text{Hyp}(N, M, L)$ , and if it is  $\ll 1$  then  $X \sim \text{Hyp}(N, M, L + 1)$  is likely to be 0 whereas  $Y \sim \text{Hyp}(N, M, L) + 1$  is at least 1. Second, it is also necessary that  $\frac{M(N-L)}{N} = M(1 - \frac{L}{N}) \gg 1$ . To see this, note that if by way of contrast  $1 - \frac{L}{N} \ll \frac{1}{M}$ , then  $X$  is concentrated at  $M$  and  $Y$  is concentrated at  $M + 1$ . Finally, to understand the hypothesis's dependence on  $e$ , suppose  $M = N/2$  and  $L$  is quite small. Then  $\text{Hyp}(N, M, L)$  is distributed very much like  $\text{Binomial}(L, e)$ ; hence we require  $L \gg e^2$  or else the extra  $+e$  in  $Y$  will dominate the standard deviation of  $\text{Binomial}(L, e)$ .

We prove Lemmas 3.4 and 3.5 in the next sections, but first we show how Theorem 3.2 follows from the lemmas.

*Proof of Theorem 3.2:* Note that we may freely assume  $k \geq 2e + 2$ , as otherwise the bound we are trying to prove exceeds 1 (assuming the constant in the  $O(\cdot)$  is large enough). Let us introduce the notation  $N = n - 1$ ,  $M(\mathbf{c}) = m(\mathbf{c}) - 1$ ,  $L = \ell - 1$ ,  $L' = \min(L, N - L)$ . Then by Lemma 3.4,

$$\begin{aligned} d_{\text{TV}}(\mathcal{S}_{\text{yes}}, \mathcal{S}_{\text{no}}) &\leq \sum_{\mathbf{c} \in \mathcal{C}: m(\mathbf{c}) \neq 0} \frac{m(\mathbf{c})}{n} \cdot \lambda_{n-1, m(\mathbf{c})-1, \ell-1}(e) \\ &= \sum_{0 \leq \frac{M(\mathbf{c})}{N} < \frac{\kappa e^2}{L'}} \frac{m(\mathbf{c})}{n} \cdot \lambda_{N, M(\mathbf{c}), L}(e) \\ &\quad + \sum_{\frac{M(\mathbf{c})}{N} \geq \frac{\kappa e^2}{L'}} \frac{m(\mathbf{c})}{n} \cdot \lambda_{N, M(\mathbf{c}), L}(e) \\ &\leq \sum_{0 \leq \frac{M(\mathbf{c})}{N} < \frac{\kappa e^2}{L'}} \frac{m(\mathbf{c})}{n} + \sum_{\frac{M(\mathbf{c})}{N} \geq \frac{\kappa e^2}{L'}} \frac{m(\mathbf{c})}{n} \cdot .01, \end{aligned}$$

where the last inequality uses Lemma 3.5. Since  $\sum_{\mathbf{c} \in \mathcal{C}} m(\mathbf{c}) = n$ , the second sum above is at most .01. Thus it remains to bound the first sum by  $|\mathcal{C}| \frac{O(e^2)}{\min(k, n-k)}$ . There are at most  $|\mathcal{C}|$  summands in this first sum, and for each we have

$$\frac{m(\mathbf{c})}{n} \leq \frac{M(\mathbf{c}) + 1}{N} \leq \frac{\kappa e^2}{L'} + \frac{1}{N} \leq \frac{2\kappa e^2}{L'}$$

by the condition of the sum.

To complete the proof, it remains to show that  $L' = \min(\ell - 1, n - \ell) \geq \Omega(\min(k, n - k))$ . When  $\ell - 1 \leq n - \ell$ , then  $L' = \ell - 1 = k - e - 1 \geq k/2$  by the fact that  $k \geq 2e + 2$ . And when  $n - \ell < \ell - 1$ , then  $L' = n - \ell = n - k + e \geq n - k$ . So  $L' \geq \frac{1}{2} \min(k, n - k)$ , as we wanted to show.  $\blacksquare$

#### D. Proof of Lemma 3.4

Let us think of the experiment  $\mathcal{S}_{\text{yes}}$  in an alternate way. We begin by choosing a first column from  $Q$  for  $\mathcal{S}_{\text{yes}}$  — call it  $C_1$ . We next decide how many additional copies of  $C_1$  to include into  $\mathcal{S}_{\text{yes}}$ . Call this quantity  $T$ . We have

$$T \mid (C_1 = \mathbf{c}) \sim \text{Hyp}(n - 1, m(\mathbf{c}) - 1, k - 1).$$

(Note that  $m(\mathbf{c}) - 1 \geq 0$  always, because  $\mathbf{c}$  won't be chosen if  $m(\mathbf{c}) = 0$ .) So far,  $\mathcal{S}_{\text{yes}}$  consists of  $T + 1$  copies of  $C_1$ . Finally, we complete the draw of  $\mathcal{S}_{\text{yes}}$  by choosing  $k - (T + 1)$  columns without replacement from " $Q \setminus C_1$ ", meaning the multiset of columns formed from  $Q$  by removing all copies of  $C_1$ .

We think of the experiment  $\mathcal{S}_{\text{no}}$  in a similar way. Again, we begin by choosing a first column  $C_1$  from  $Q$  for  $\mathcal{S}_{\text{no}}$ . We next determine how many additional copies of  $C_1$  there will be from among the remaining  $\ell - 1$  choices. Calling this quantity  $U$ , we have

$$U \mid (C_1 = \mathbf{c}) \sim \text{Hyp}(n - 1, m(\mathbf{c}) - 1, \ell - 1).$$

Recall, however, that in  $\mathcal{S}_{\text{no}}$ , we include an additional  $e$  copies of  $C_1$  into  $\mathcal{S}_{\text{no}}$ . Hence  $\mathcal{S}_{\text{no}}$  ends up with  $U + e + 1$  copies of  $C_1$ . Finally, we complete  $\mathcal{S}_{\text{no}}$  by adding  $\ell - (U + 1)$  columns drawn without replacement from  $Q \setminus C_1$ .

Let  $V = U + e$ . We claim that by coupling the random variables  $T \mid (C_1 = \mathbf{c})$  and  $V \mid (C_1 = \mathbf{c})$ , we couple  $\mathcal{S}_{\text{yes}}$  and  $\mathcal{S}_{\text{no}}$ . This follows immediately from the two descriptions, as then  $T + 1 = V + 1 = U + e + 1$ , and  $k - (T + 1) = \ell + e - (V + 1) = \ell - (U + 1)$ . Hence

$$\begin{aligned} d_{\text{TV}}(\mathcal{S}_{\text{yes}}, \mathcal{S}_{\text{no}}) &\leq \sum_{\mathbf{c} \in \mathcal{C}} \Pr[C_1 = \mathbf{c}] \cdot \\ &\quad d_{\text{TV}}(T \mid (C_1 = \mathbf{c}), V \mid (C_1 = \mathbf{c})). \end{aligned}$$

On one hand,  $\Pr[C_1 = \mathbf{c}]$  is simply  $\frac{m(\mathbf{c})}{n}$ . On the other hand, we have

$$\begin{aligned} T \mid (C_1 = \mathbf{c}) &\sim \text{Hyp}(n - 1, m(\mathbf{c}) - 1, \ell + e - 1), \\ V \mid (C_1 = \mathbf{c}) &\sim \text{Hyp}(n - 1, m(\mathbf{c}) - 1, \ell - 1) + e. \end{aligned}$$

So by definition,  $d_{\text{TV}}(T \mid (C_1 = \mathbf{c}), V \mid (C_1 = \mathbf{c})) = \lambda_{n-1, m(\mathbf{c})-1, \ell-1}(e)$ , and hence

$$d_{\text{TV}}(\mathcal{S}_{\text{yes}}, \mathcal{S}_{\text{no}}) \leq \sum_{\mathbf{c} \in \mathcal{C}: m(\mathbf{c}) \neq 0} \frac{m(\mathbf{c})}{n} \cdot \lambda_{n-1, m(\mathbf{c})-1, \ell-1}(e),$$

as claimed.

### E. Proof of Lemma 3.5

Recall that  $L' = \min(L, N - L)$ ,

$$\frac{ML'}{N} \geq \kappa e^2, \quad (3)$$

and our goal is to bound  $\lambda_{N,M,L}(e) = d_{\text{TV}}(X, Y) \leq .01$ , where  $X \sim \text{Hyp}(N, M, L+e)$  and  $Y \sim \text{Hyp}(N, M, L)+e$ .

We begin by coupling  $X$  and  $Y$ , as follows. Imagine drawing balls randomly and without replacement from an urn containing  $N$  balls,  $M$  of which are white. We draw  $L+e$  balls from the urn. We let  $X$  be the number of white balls among all balls drawn; we let  $Y$  be the number of white balls among the first  $L$  balls drawn, plus  $e$ . Note that  $X \leq Y$  always under this coupling.

Let us now compare the probability mass functions of  $X$  and  $Y$ . The integers  $u < e$  can be in  $X$ 's range but not  $Y$ 's; the integers  $u > \min(M, L+e)$  can be in  $Y$ 's range but not  $X$ 's. The remaining integers are in the range of both  $X$  and  $Y$ , and we have

$$\begin{aligned} \frac{\Pr[X = u]}{\Pr[Y = u]} &= \frac{\binom{M}{u} \binom{N-M}{L+e-u}}{\binom{N}{L+e}} \bigg/ \frac{\binom{M}{u-e} \binom{N-M}{L+e-u}}{\binom{N}{L}} \\ &= \frac{\binom{M}{u}}{\binom{M}{u-e}} \cdot \frac{\binom{N}{L}}{\binom{N}{L+e}} \\ &= \frac{(M-u+e)(M-u+e-1)\cdots(M-u+1)}{u(u-1)\cdots(u-e+1)} \cdot \frac{\binom{N}{L}}{\binom{N}{L+e}}. \end{aligned}$$

Evidently (and unsurprisingly), this ratio is a decreasing function of  $u$ . Letting  $t$  be the largest integer for which the ratio is at least 1, we conclude that

$$\Pr[X = u] \geq \Pr[Y = u] \text{ iff } u \leq t.$$

It follows immediately that

$$d_{\text{TV}}(X, Y) = \Pr[X \leq t] - \Pr[Y \leq t].$$

But by our coupling,

$$\begin{aligned} \Pr[X \leq t] - \Pr[Y \leq t] &= \Pr[X \leq t \cap Y > t] \\ &\quad - \Pr[X > t \cap Y \leq t] \\ &= \Pr[X \leq t \cap Y > t], \end{aligned}$$

since  $X \leq Y$  always. Our goal, then, is to bound

$$d_{\text{TV}}(X, Y) = \Pr[X \leq t \cap Y > t]. \quad (4)$$

We will in fact prove something slightly stronger: we will show that for *any* value of  $t$ , the right-hand side of (4) is small.

To analyze (4) we recall the ball and urn process defining  $X$  and  $Y$ . Having drawn  $L+e$  balls, let  $W$  be the number of white balls among the *last*  $e$  balls drawn, and let  $Z$  be the number of white balls among the first  $L$ . Thus  $X = W + Z$  and  $Y = e + Z$ . As a first observation, we may note that

if  $W = e$  then  $X = Y$  and hence the event in (4) does not occur. I.e.,

$$d_{\text{TV}}(X, Y) \leq \Pr[W \neq e] \leq e(1 - \frac{M}{N}), \quad (5)$$

where we used a union bound over each of the last  $e$  balls being non-white. Now by (3),

$$e \leq \sqrt{\frac{1}{\kappa} \frac{ML'}{N}} \leq \sqrt{\frac{L'}{\kappa}} \leq .001\sqrt{N}, \quad (6)$$

if we assume  $\kappa$  large enough. It follows that we may additionally assume

$$M \leq N - .01\sqrt{N} \quad \Leftrightarrow \quad 1 - \frac{M}{N} \geq \frac{.01}{\sqrt{N}} \quad (7)$$

because otherwise the bound in (5) is at most  $.001\sqrt{N} \cdot \frac{.01}{\sqrt{N}} = .00001$ , which establishes the theorem with room to spare. We also use this opportunity to mention that

$$M \geq 2e, \quad L' \geq 2e, \quad (\text{and hence certainly } N \geq 2e) \quad (8)$$

follow easily from (3).

We next give a more refined upper bound on (4). By conditioning on  $W$  we have

$$\begin{aligned} d_{\text{TV}}(X, Y) &= \Pr[X \leq t \cap Y > t] \\ &= \sum_{i=0}^{e-1} \Pr[W = i] \Pr[t - e < Z \leq t - i \mid W = i]. \end{aligned}$$

Now  $Z \mid (W = i)$  has distribution  $\text{Hyp}(N - e, M - i, L)$  (and note that  $M - i \geq M - e \geq 0$  by (8)). Let us write  $\sigma^2 = L(1 - \frac{L}{N-e}) \frac{M-i}{N-e} (1 - \frac{M-i}{N-e})$ . Applying Corollary 2.1 and a union bound we get

$$\begin{aligned} d_{\text{TV}}(X, Y) &\leq \sum_{i=0}^{e-1} \Pr[W = i] \cdot (e - i) \frac{C}{\sigma} \\ &\leq \max_{0 \leq i < e} \left\{ \frac{C}{\sigma} \right\} \cdot \sum_{i=0}^{e-1} \Pr[W = i] (e - i) \\ &= \max_{0 \leq i < e} \left\{ \frac{C}{\sigma} \right\} \cdot \mathbf{E}[e - W]. \end{aligned}$$

We have  $W \sim \text{Hyp}(N, M, e)$ , and thus  $\mathbf{E}[e - W] = e(1 - \frac{M}{N})$ . And by definition,

$$\begin{aligned} \max_{0 \leq i < e} \left\{ \frac{C}{\sigma} \right\} &= \max_{0 \leq i < e} \left\{ \frac{C}{\sqrt{L(1 - \frac{L}{N-e})(\frac{M-i}{N-e})(1 - \frac{M-i}{N-e})}} \right\} \\ &\leq \frac{C}{\sqrt{L(1 - \frac{L}{N-e})(\frac{M-e}{N-e})(1 - \frac{M}{N-e})}}. \end{aligned}$$

Thus we have established

$$d_{\text{TV}}(X, Y) \leq \frac{Ce}{\sqrt{L(1 - \frac{L}{N-e})}} \cdot \frac{1 - \frac{M}{N}}{\sqrt{1 - \frac{M}{N-e}}} \cdot \frac{1}{\sqrt{\frac{M-e}{N-e}}}. \quad (9)$$

We will bound the three fractions in (9) one at a time. We begin with the middle one. Note first that

$$\frac{d}{dM} \left( \frac{1 - \frac{M}{N}}{\sqrt{1 - \frac{M}{N-e}}} \right) = -\frac{N - 2e - M}{2N\sqrt{1 - \frac{M}{N-e}}(N - e - M)}.$$

By combining (6) and (7) we get  $M \leq N - 10e < N - 2e$ . Hence the derivative above is always negative, implying that  $(1 - \frac{M}{N})/\sqrt{1 - \frac{M}{N-e}}$  is a decreasing function of  $M$  on  $M$ 's range. Hence we may upper-bound this fraction by taking  $M = 0$ , giving an upper bound of 1. Substituting this into (9) gives

$$d_{\text{TV}}(X, Y) \leq \frac{Ce}{\sqrt{L(1 - \frac{L}{N-e})}} \cdot \frac{1}{\sqrt{\frac{M-e}{N-e}}}. \quad (10)$$

We next examine the fraction on the right. It is at most

$$\frac{1}{\sqrt{\frac{M-e}{N}}} \leq \frac{1}{\sqrt{\frac{M/2}{N}}} = \sqrt{\frac{2N}{M}},$$

where we used (8). By virtue of (3), we can upper-bound this by  $\sqrt{\frac{2}{\kappa}} \cdot \frac{\sqrt{L'}}{e}$ . Substituting this upper bound into (10) yields

$$\begin{aligned} d_{\text{TV}}(X, Y) &\leq C\sqrt{\frac{2}{\kappa}} \cdot \sqrt{\frac{L'}{L(1 - \frac{L}{N-e})}} \\ &\leq .001\sqrt{\frac{L'}{L(1 - \frac{L}{N-e})}}, \end{aligned} \quad (11)$$

assuming  $\kappa$  is sufficiently large compared with  $C$ .

Finally, we split into two cases, depending on whether  $L \leq N/2$ . If indeed  $L \leq N/2$ , then  $L' = L$  and we have

$$.001\sqrt{\frac{L'}{L(1 - \frac{L}{N-e})}} = \frac{.001}{\sqrt{1 - \frac{L}{N-e}}} \leq \frac{.001}{\sqrt{1 - \frac{N/2}{N-e}}}.$$

But  $N - e \geq N - .001\sqrt{N} \geq (2/3)N$  (using (6) and  $N \geq 2$  from (8)), so we upper-bound

$$d_{\text{TV}}(X, Y) \leq \frac{.001}{\sqrt{1 - \frac{N/2}{(2/3)N}}} = .002 \leq .01,$$

as needed. The second case is that  $L \geq N/2$ , in which case  $L = N - L'$  and the bound in (11) is

$$\begin{aligned} .001\sqrt{\frac{L'}{(N - L')(1 - \frac{N-L'}{N-e})}} &= .001\sqrt{\frac{L'}{(N - L')\frac{L'-e}{N-e}}} \\ &= .001\sqrt{\frac{L'}{L'-e}} \sqrt{\frac{N-e}{N-L'}}. \end{aligned} \quad (12)$$

But using (8),

$$\sqrt{\frac{L'}{L'-e}} \leq \sqrt{\frac{L'}{L'/2}} = \sqrt{2},$$

and using  $L' \leq N/2$ ,

$$\sqrt{\frac{N-e}{N-L'}} \leq \sqrt{\frac{N}{N-L'}} \leq \sqrt{\frac{N}{N/2}} = \sqrt{2}.$$

Hence the upper bound (12) on  $d_{\text{TV}}(X, Y)$  is at most  $.001\sqrt{2}\sqrt{2} < .01$ , as needed.

This completes the proof of Lemma 3.5.

#### IV. MAJORITY FUNCTIONS

Recall that  $\text{Maj}_k : \{-1, 1\}^n \rightarrow \{-1, 1\}$  is defined by  $\text{Maj}_k(x) = \text{sgn}(\sum_{i=1}^k x_i)$ , where we define  $\text{sgn}(0)$  (arbitrarily) to be 1. We sometimes abuse notation by thinking of  $\text{Maj}_k$  also as a function on  $\{-1, 1\}^k$ . In this section we prove our lower bound for testing  $\text{Maj}_k$ -isomorphism, Theorem 1.5, which we restate here for convenience:

**Theorem 1.5 (Restated).** *There is a universal  $\epsilon_0 > 0$  such that whenever  $k \leq (3/4)n$  (and is divisible by 3), any non-adaptive algorithm for  $\epsilon_0$ -testing  $\text{Maj}_k$ -isomorphism must make at least  $\Omega(k^{1/12})$  queries.*

Note that the result of Theorem 1.5 cannot be handled by Theorem 1.3 since the  $\text{Maj}_k$  function is  $o(1)$ -close to being a  $(k - e)$ -junta whenever  $e = o(k)$ . Specifically, the full version of this paper establishes the following proposition, which is very similar to the problem of computing the noise sensitivity of majority [16].

**Proposition 4.1.** *For  $0 < e < k/2$ , the  $\text{Maj}_k$  function is  $\epsilon$ -close to the  $\text{Maj}_{k-e}$  function, where  $\epsilon = 6(e/k)^{1/3}$ .*

Our proof of Theorem 1.5 reuses much of the framework introduced in the previous section in our testing lower bound for general  $g$ . As before, our goal is to construct probability distributions  $\mathcal{F}_{\text{yes}}$  and  $\mathcal{F}_{\text{no}}$  over functions isomorphic to  $\text{Maj}_k$  and functions  $\epsilon_0$ -far from being isomorphic to  $\text{Maj}_k$  (respectively) such that any deterministic non-adaptive testing algorithm making  $o(k^{1/12})$  queries cannot distinguish with probability at least  $1/3$  between functions drawn from  $\mathcal{F}_{\text{yes}}$  or from  $\mathcal{F}_{\text{no}}$ .

We define  $\mathcal{F}_{\text{yes}}$  just as we did in Section III: a function  $f_{\text{yes}} \sim \mathcal{F}_{\text{yes}}$  is obtained by choosing  $j_1, \dots, j_k$  randomly and without replacement from  $[n]$  and defining  $f_{\text{yes}}(x) = \text{Maj}_k(x_{j_1}, \dots, x_{j_k})$ . The definition  $\mathcal{F}_{\text{no}}$ , however, is very different from the definition we used in that section to ensure that it is supported on functions far from  $\text{Maj}_k$ .

To define  $\mathcal{F}_{\text{no}}$ , we first introduce a certain *weighted* majority function  $\text{WgtMaj}_k$  on  $\frac{4}{3}k$  bits (note that  $\frac{4}{3}k \leq n$ ):

$$\text{WgtMaj}_k(x_1, \dots, x_{\frac{4}{3}k}) =$$

$$\text{sgn}\left(\sum_{i=1}^{k/3} \left(\frac{1}{2}x_{4i-3} + \frac{1}{2}x_{4i-2} + \frac{1}{2}x_{4i-1} + \frac{3}{2}x_{4i}\right)\right). \quad (13)$$



I.e.,  $\text{WgtMaj}_k$  gives  $k$  variables weight  $\frac{1}{2}$  and  $k/3$  variables weight  $\frac{3}{2}$ . This weight pattern is chosen very carefully; see the proof of Lemma 4.7 below. We then define  $f_{\text{no}} \sim \mathcal{F}_{\text{no}}$  is obtained by choosing  $j_1, \dots, j_{\frac{4}{3}k}$  randomly and without replacement from  $[n]$  and taking  $f_{\text{no}}(x) = \text{WgtMaj}_k(x_{j_1}, \dots, x_{j_{\frac{4}{3}k}})$ .

The following proposition, whose proof is included in the full version of this paper, implies that functions in  $\mathcal{F}_{\text{no}}$  are  $\epsilon$ -far from  $\text{Maj}_k$  functions for large enough values of  $\epsilon$ :

**Proposition 4.2.** *There exist universal constants  $\epsilon_0 > 0$  and  $k_0 \in \mathbb{N}$  such that when  $k \geq k_0$ , every function  $f_{\text{no}}$  in the support of  $\mathcal{F}_{\text{no}}$  is  $\epsilon_0$ -far from being a  $k$ -junta.*

Note that we may always assume  $k \geq k_0$  as otherwise Theorem 1.5 is trivial. To complete the proof of Theorem 1.5, it suffices to prove the following:

**Theorem 4.3.** *Let  $\mathcal{T}$  be any deterministic non-adaptive  $q$ -query algorithm for testing isomorphism to  $\text{Maj}_k$ . Then*

$$\left| \Pr_{f_{\text{yes}} \sim \mathcal{F}_{\text{yes}}} [\mathcal{T} \text{ acc. } f_{\text{yes}}] - \Pr_{f_{\text{no}} \sim \mathcal{F}_{\text{no}}} [\mathcal{T} \text{ acc. } f_{\text{no}}] \right| \leq O\left(\frac{q^{3/2}}{k^{1/8}}\right).$$

To prove Theorem 4.3, we continue to recall the framework developed in Section III. Given a deterministic  $q$ -query tester  $\mathcal{T}$ , we arrange its  $q$  queries  $x^1, \dots, x^q \in \{-1, 1\}^n$  into a  $q \times n$  query matrix  $Q$ . We again think of two processes  $\mathcal{S}_{\text{yes}}$  and  $\mathcal{S}_{\text{no}}$  for generating argument multisets  $S_{\text{yes}}$  and  $S_{\text{no}}$ . However in the present case we simply have that  $\mathcal{S}_{\text{yes}}$  chooses  $k$  columns at random from  $Q$  without replacement, and  $\mathcal{S}_{\text{no}}$  chooses  $\frac{4}{3}k$  columns at random from  $Q$  without replacement. Again, we imagine that the argument multisets are randomly ordered to form argument matrices:  $A_{\text{yes}}$  which is  $q \times k$ , and  $A_{\text{no}}$  which is  $q \times \frac{4}{3}k$ . Finally, we obtain the response vector random variable  $R_{\text{yes}} \in \{-1, 1\}^q$  by applying  $\text{Maj}_k$  to  $A_{\text{yes}}$  row-wise, and the response vector  $R_{\text{no}} \in \{-1, 1\}^q$  by applying  $\text{WgtMaj}_k$  to  $A_{\text{no}}$  row-wise. It is clear that this distribution on  $R_{\text{yes}}$  is equivalent to the one given by drawing  $f_{\text{yes}} \sim \mathcal{F}_{\text{yes}}$  and letting  $R_{\text{yes}} = \langle f_{\text{yes}}(x^1), \dots, f_{\text{yes}}(x^q) \rangle$ . The analogous statement is true for  $R_{\text{no}}$ . Hence we can prove Theorem 4.3 by showing

$$d_{\text{TV}}(R_{\text{yes}}, R_{\text{no}}) \leq O\left(\frac{q^{3/2}}{k^{1/8}}\right). \quad (14)$$

We now come to the main difference between our  $\text{Maj}_k$  lower bound and the general lower bound from Section III. Obviously, we cannot proceed as in Section III by bounding the total variation distance between  $S_{\text{yes}}$  and  $S_{\text{no}}$ : this total variation distance is 1, since  $\mathcal{S}_{\text{yes}}$  and  $\mathcal{S}_{\text{no}}$  have disjoint support! (Specifically, the multiset  $S_{\text{yes}}$  has cardinality  $k$  whereas  $S_{\text{no}}$  has cardinality  $\frac{4}{3}k$ .) Instead, we exploit the fact that applying  $\text{Maj}_k$  or  $\text{WgtMaj}_k$  involves *adding up* the columns in the argument matrix (in  $\text{WgtMaj}_k$ 's case, with certain weights), and this addition ‘‘loses a lot of information’’.

More precisely, suppose that we write  $X_1, \dots, X_k$  for the (randomly chosen) columns in argument matrix  $A_{\text{yes}}$  and let

$$S = X_1 + \dots + X_k.$$

Then the response vector  $R_{\text{yes}}$  is given by taking the  $\text{sgn}$  of each entry of  $S$ ; i.e., it is determined by the orthant of  $\mathbb{R}^q$  in which  $S$  lies. Similarly, if we write  $Y_1, \dots, Y_{\frac{4}{3}k}$  for the columns in argument matrix  $A_{\text{no}}$ , and let

$$T = \frac{1}{2}Y_1 + \frac{1}{2}Y_2 + \frac{1}{2}Y_3 + \frac{3}{2}Y_4 + \dots \\ + \frac{1}{2}Y_{\frac{4}{3}k-3} + \frac{1}{2}Y_{\frac{4}{3}k-2} + \frac{1}{2}Y_{\frac{4}{3}k-1} + \frac{3}{2}Y_{\frac{4}{3}k},$$

then  $R_{\text{no}}$  is determined by the orthant in which  $T$  lies. Hence we can establish (14) and thus Theorem 1.5 by proving the following:

**Theorem 4.4.** *Let  $S$  and  $T$  be defined as above. Then for any union  $\mathcal{O}$  of orthants in  $\mathbb{R}^d$ ,*

$$|\Pr[S \in \mathcal{O}] - \Pr[T \in \mathcal{O}]| \leq O(q^{3/2}/k^{1/8}).$$

#### A. Multidimensional invariance

To prove Theorem 4.4 we will use recently developed invariance principle tools [14], [13], [9]. In particular, we quote the following multidimensional results which essentially appear in [13], [9].

**Lemma 4.5.** *(Essentially Theorem 4.1 in [13]; cf. [9].) Let  $\psi : \mathbb{R}^q \rightarrow \mathbb{R}$  be a thrice continuously differentiable function with uniformly bounded third partial derivatives:  $|\psi^{(J)}| \leq \beta$  for all multi-indices  $J = (j_1, \dots, j_q)$  with  $|J| = j_1 + \dots + j_q = 3$ . Let  $S = S_1 + \dots + S_m$ , where the  $S_i$ 's are independent  $\mathbb{R}^q$ -valued random variables, and let  $T = T_1 + \dots + T_m$  similarly. Assume that for each  $i \in [m]$ ,  $S_i$  and  $T_i$  have matching means and covariance matrices:  $\mathbf{E}[S_i] = \mathbf{E}[T_i]$  and  $\text{Cov}[S_i] = \text{Cov}[T_i]$ . Then*

$$|\mathbf{E}[\psi(S)] - \mathbf{E}[\psi(T)]| \leq O(\beta) \cdot \sum_{i=1}^m \sum_{|J|=3} \left( \mathbf{E}[|S_i^J|] + \mathbf{E}[|T_i^J|] \right),$$

where  $U^J$  denotes  $U_1^{j_1} \dots U_q^{j_q}$  when  $U \in \mathbb{R}^q$ .

**Lemma 4.6.** *(Essentially appears in [9].) Let  $\mathcal{O}$  be any union of orthants in  $\mathbb{R}^q$  and let  $S, T$  be any  $\mathbb{R}^q$ -valued random variables. Let  $r > 0$ . Then there is a certain smooth function  $\psi$  satisfying  $|\psi^{(J)}| \leq O(1/r^3)$  for all multi-indices  $J$  with  $|J| = 3$  and such that*

$$|\Pr[S \in \mathcal{O}] - \Pr[T \in \mathcal{O}]| \leq \Pr[S \in W_r] + |\mathbf{E}[\psi(S)] - \mathbf{E}[\psi(T)]|,$$

where  $W_r = \{x \in \mathbb{R}^q : |x_i| \leq r/2 \text{ for some } i \in [q]\}$ .

At first, it does not appear as though these tools are of any help to us, because Lemma 4.5 very crucially uses the fact that the random vectors being summed are independent. Whereas, in our Theorem 4.4 the random vectors  $X_1, \dots, X_k$  are certainly not independent, being drawn randomly *without replacement* from the fixed population  $Q$ . The same goes for  $Y_1, \dots, Y_{\frac{4}{3}k}$ . Nevertheless, we can still reduce to Lemma 4.5 using a trick: finding random vectors which are *conditionally independent*.

### B. How to handle drawing without replacement

Let us recap the scenario in Theorem 4.4. We have a fixed multiset  $Q$  of  $n$  columns (vectors) from  $\{-1, 1\}^q$ . We draw  $k$  columns randomly from  $Q$  without replacement, yielding the vector-valued random variables  $X_1, \dots, X_k$ ; we also define

$$S = X_1 + \dots + X_k.$$

Similarly, the vector-valued random variables  $Y_1, \dots, Y_{\frac{4}{3}k}$  are drawn randomly from  $Q$  without replacement, and

$$T = \frac{1}{2}Y_1 + \frac{1}{2}Y_2 + \frac{1}{2}Y_3 + \frac{3}{2}Y_4 + \dots \\ + \frac{1}{2}Y_{\frac{4}{3}k-3} + \frac{1}{2}Y_{\frac{4}{3}k-2} + \frac{1}{2}Y_{\frac{4}{3}k-1} + \frac{3}{2}Y_{\frac{4}{3}k}.$$

To introduce conditional independence, we reimagine how the  $X_i$ 's and  $Y_i$ 's are drawn. Specifically, we can couple the drawing of these random vectors as follows:

- 1) Define  $m = k/3$  (an integer).
- 2) Randomly partition the columns of  $Q$  into  $m$  parts  $Q_1, \dots, Q_m$ , each of cardinality 4, along with a leftover set of cardinality  $n - 4m \geq 0$ .
- 3) Independently for each  $i \in [m]$ , choose  $X_{3i-2}, X_{3i-1}, X_{3i}$  randomly without replacement from  $Q_i$ . Define also  $S_i = X_{3i-2} + X_{3i-1} + X_{3i}$ .
- 4) Independently for each  $i \in [m]$ , choose  $Y_{4i-3}, Y_{4i-2}, Y_{4i-1}, Y_{4i}$ , randomly without replacement from  $Q_i$  (i.e., choose them by randomly ordering the vectors in  $Q_i$ ). Define also  $T_i = \frac{1}{2}Y_{4i-3} + \frac{1}{2}Y_{4i-2} + \frac{1}{2}Y_{4i-1} + \frac{3}{2}Y_{4i}$ .

It is easy to see that this coupling gives the correct marginal distributions on  $X_1, \dots, X_k$  and  $Y_1, \dots, Y_{\frac{4}{3}k}$ . We also have  $S = S_1 + \dots + S_m$  and  $T = T_1 + \dots + T_m$ . And crucially,  $S_1, \dots, S_m$  are independent *conditioned on any choice of the partition*  $(Q_1, \dots, Q_m)$ , and similarly for  $T_1, \dots, T_m$ . The following lemma will allow us to apply Lemma 4.5; it also explains the choice of the weight pattern  $(\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{3}{2})$ :

**Lemma 4.7.** *Conditioned on any choice of the partition  $(Q_1, \dots, Q_m)$ , we have  $\mathbf{E}[S_i] = \mathbf{E}[T_i]$  and  $\mathbf{Cov}[S_i] = \mathbf{Cov}[T_i]$  for each  $i \in [m]$ .*

Let  $W_r = \{x \in \mathbb{R}^d : \exists j \in [d] \text{ s.t. } |x_j| \leq r\}$  represent the region around the orthant boundaries. The following Lemma gives an upper bound on the probability that our random vector  $S$  lands near any orthant boundary:

**Lemma 4.8.** *For  $r \geq 1$  it holds that*

$$\Pr[S \in W_r] \leq O\left(\frac{qr}{\sqrt{m}}\right).$$

We present the proofs of Lemmas 4.7 and 4.8 below. First, let us now combine all these results to complete the proof of Theorem 4.4 and hence Theorem 1.5:

*Proof of Theorem 4.4:* Let us first condition on a particular partition  $(Q_1, \dots, Q_m)$ . Having done so,  $S_1, \dots, S_m$  become independent, as do  $T_1, \dots, T_m$ . By Lemma 4.7, we

may apply Lemma 4.5. Doing so with the function  $\psi$  from Lemma 4.6 (with  $r \geq 1$  to be chosen later) yields

$$|\mathbf{E}[\psi(S)] - \mathbf{E}[\psi(T)]| \leq O\left(\frac{1}{r^3}\right) \cdot \sum_{i=1}^m \sum_{|J|=3} \left(\mathbf{E}[|S_i^J|] + \mathbf{E}[|T_i^J|]\right). \quad (15)$$

We emphasize that (15) is conditional on a particular  $(Q_1, \dots, Q_m)$ . However, note that we can bound the quantities  $\mathbf{E}[|S_i^J|]$  and  $\mathbf{E}[|T_i^J|]$  uniformly in  $(Q_1, \dots, Q_m)$ : Each coordinate of  $S_i$  is at most  $1 + 1 + 1 = 3$  in absolute value, and similarly each coordinate of  $T_i$  is at most  $\frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{3}{2} = 3$  in absolute value. Hence each expectation is at most 27, and we can therefore upper-bound the right-hand side of (15) by  $O(mq^3/r^3)$ , since there are at most  $q^3$  many  $J$ 's.

Now averaging over the choice of partition  $(Q_1, \dots, Q_m)$ , the triangle inequality implies

$$|\mathbf{E}[\psi(S)] - \mathbf{E}[\psi(T)]| \leq O(mq^3/r^3).$$

Here the expectation is over the whole definition of  $S$  and  $T$ . Substituting this into Lemma 4.6 gives

$$|\Pr[S \in \mathcal{O}] - \Pr[T \in \mathcal{O}]| \leq \Pr[S \in W_r] + O(mq^3/r^3).$$

Applying Lemma 4.8 we can bound this by  $O(qr/\sqrt{m}) + O(mq^3/r^3)$ . We optimize by taking  $r = m^{3/8}d^{1/2}$ , yielding a final upper bound of  $O(q^{3/2}/m^{1/8}) = O(q^{3/2}/k^{1/8})$  and completing the proof. ■

### C. Proof of Lemma 4.7

It suffices to check the claim for  $i = 1$ . Let  $\mu \in \mathbb{R}^q$  denote the average of the four vectors in  $Q_1$ . Then

$$\mathbf{E}[S_1] = \mathbf{E}[X_1 + X_2 + X_3] = 3\mathbf{E}[X_1] = 3\mu,$$

using the fact that  $X_1, X_2, X_3$  are identically distributed, and similarly

$$\mathbf{E}[T_1] = \mathbf{E}\left[\frac{1}{2}Y_1 + \frac{1}{2}Y_2 + \frac{1}{2}Y_3 + \frac{3}{2}Y_4\right] \\ = \left(\frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{3}{2}\right)\mathbf{E}[Y_1] = 3\mu,$$

verifying the first claim.

We now verify that  $\mathbf{Cov}[S_i] = \mathbf{Cov}[T_i]$ . Fix  $j, j' \in [q]$  and let us write  $x_1, x_2, x_3$  for the  $j$ -coordinate of  $X_1, X_2, X_3$  respectively, and write  $x'_1, x'_2, x'_3$  for the  $j'$ -coordinate of  $X_1, X_2, X_3$ . Then

$$\mathbf{E}[(S_1)_j(S_1)_{j'}] = \mathbf{E}[(x_1 + x_2 + x_3)(x'_1 + x'_2 + x'_3)] \\ = 3\mathbf{E}[x_1x'_1] + 6\mathbf{E}[x_1x'_2],$$

where we use the facts that  $(x_1, x_1)$  has the same distribution as  $(x_2, x_2)$  and  $(x_3, x'_3)$ , and that  $(x_1, x'_2)$  has the same distribution as  $(x_\ell, x'_{\ell'})$  for any  $\ell \neq \ell' \in \{1, 2, 3\}$ . So

$$\mathbf{Cov}[S_1]_{j,j'} = 3\mathbf{E}[x_1x'_1] + 6\mathbf{E}[x_1x'_2] - \mu_j\mu_{j'}.$$

Similarly, we define  $y_1, \dots, y_4$  and  $y'_1, \dots, y'_4$  as the  $j$ th and  $j'$ th coordinates of  $Y_1, \dots, Y_4$ , respectively. Then

$$\begin{aligned} \mathbf{E}[(T_1)_j(T_1)_{j'}] &= \mathbf{E}[(\tfrac{1}{2}y_1 + \tfrac{1}{2}y_2 + \tfrac{1}{2}y_3 + \tfrac{3}{2}y_4) \cdot \\ &\quad (\tfrac{1}{2}y'_1 + \tfrac{1}{2}y'_2 + \tfrac{1}{2}y'_3 + \tfrac{3}{2}y'_4)] \\ &= (3 \cdot (\tfrac{1}{2})^2 + (\tfrac{3}{2})^2) \mathbf{E}[y_1 y'_1] + \\ &\quad (6 \cdot (\tfrac{1}{2})^2 + 6 \cdot \tfrac{1}{2} \cdot \tfrac{3}{2}) \mathbf{E}[y_1 y'_2] \\ &= 3 \mathbf{E}[y_1 y'_1] + 6 \mathbf{E}[y_1 y'_2], \end{aligned}$$

where we use the facts that  $(y_1, y'_1)$  has the same distribution as  $(y_\ell, y'_\ell)$  for any  $\ell = 2, 3, 4$  and  $(y_1, y'_2)$  has the same distribution as  $(y_\ell, y'_\ell)$  for any  $\ell \neq \ell' \in \{1, 2, 3, 4\}$ . So

$$\mathbf{Cov}[T_1]_{j,j'} = 3 \mathbf{E}[y_1 y'_1] + 6 \mathbf{E}[y_1 y'_2] - \mu_j \mu_{j'}.$$

Noting that  $(X_1, X_2)$  and  $(Y_1, Y_2)$  have the same distribution, we get that  $\mathbf{Cov}[S_1] = \mathbf{Cov}[T_1]$ .

#### D. Proof of Lemma 4.8

By union-bounding over the  $q$  coordinates, Lemma 4.8 reduces to proving the following statement:

**Lemma 4.9.** *Let  $r \geq 1$ . Suppose we fix any query row  $(x_1, \dots, x_n) \in \{-1, 1\}^n$  from  $Q$  and form the random variable*

$$s = x_{i_1} + \dots + x_{i_{3m}},$$

where the sequence  $i_1, \dots, i_{3m}$  is drawn randomly without replacement from  $[n]$ . Then

$$\Pr[|s| \leq r/2] \leq O(r/\sqrt{m}).$$

*Proof:* Let us recall that  $k = 3m \leq (3/4)n$ . Let  $u$  denote the number of 1's among  $x_1, \dots, x_n$ . The statement to be proved is precisely equivalent to the following: Let  $Z \sim \text{Hyp}(n, u, k)$ . Then

$$\Pr[Z \in [k/2 - r/4, k/2 + r/4]] \leq O(r/\sqrt{k}). \quad (16)$$

We divide into two cases.

*Case 1:*  $1/4 \leq u/n \leq 3/4$ : In this case we use Corollary 2.1 and a union bound over the at most  $r/2 + 1$  integers in the range  $[k/2 - r/4, k/2 + r/4]$  to deduce

$$\Pr[Z \in [k/2 - r/4, k/2 + r/4]] \leq O(r)/\sigma_{n,u,k},$$

where  $\sigma_{n,u,k} = \sqrt{k(1 - k/n)(u/n)(1 - u/n)}$ . We have  $1 - k/n \geq 1/4$  and also  $u/n, 1 - u/n \geq 1/4$ . Thus  $\sigma_{n,u,k} = \Omega(\sqrt{k})$ , establishing (16).

*Case 2:*  $u/n \notin [1/4, 3/4]$ : By symmetry, it suffices to treat just one of the cases  $u/n < 1/4$  or  $u/n > 3/4$ ; say, the former. In this case we have  $\mathbf{E}[Z] = k(u/n) \leq k/4$ , and we have  $\mathbf{Var}[Z] = ku/n(1 - u/n)(1 - k/n) \leq k$ . Finally, we may assume that  $r \leq k/2$ , as otherwise (16) is trivial. Thus by Chebyshev's Inequality,

$$\Pr[Z \geq k/2 - r/4] \leq \Pr[Z \geq (3/8)k] \leq \frac{k}{(k/8)^2} = O(\tfrac{1}{k}),$$

establishing (16) with room to spare.  $\blacksquare$

## V. DISCUSSION

We conclude this work by discussing what we feel are promising directions towards closing the basic research problem of characterizing the functions  $g$  for which  $g$ -isomorphism is testable.

It is possible that Fischer et al. [7]'s positive result, that testing  $g$ -isomorphism is easy when  $g$  is an  $O(1)$ -junta is mostly best possible. We pose the following question:

**Question:** Suppose  $g : \{0, 1\}^n \rightarrow \{0, 1\}$  is an  $(n - \ell)$ -junta which is  $\epsilon$ -far from being an  $\ell$ -junta. Is it true that  $\epsilon$ -testing  $g$ -isomorphism requires  $\omega_\ell(1)$  queries?<sup>2</sup>

We are not quite bold enough to conjecture that this is true, but we do not know any  $g$  which rules it out. Proving the result seems like it might be difficult, but we believe the problem is approachable for the special case when  $g$  is a *symmetric*  $k$ -junta. Our Theorem 1.4 establishes the result for the simplest symmetric function,  $\text{Maj}_k$ . It is not too hard to extend our methods to deal with similar symmetric functions; for example, we are able to show (proof omitted) that testing  $g$ -isomorphism is hard for a function  $g : \{0, 1\}^n \rightarrow \{0, 1\}$  such as

$$g(x) = \begin{cases} 1 & \text{if } k/2 - \sqrt{k} \leq \sum_{i=1}^k x_i \leq k/2 + \sqrt{k}, \\ 0 & \text{else.} \end{cases}$$

Roughly speaking, we can handle this case because it has only a constant number of ‘‘jumps’’ (just 2, in fact) between 0 and 1 on the main range of  $\sum_{i=1}^k x_i$ , namely  $k/2 \pm O(\sqrt{k})$ . If, on the other hand,  $g$  is a symmetric  $k$ -junta with ‘‘many’’ jumps between 0 and 1 on the main range of  $\sum_{i=1}^k x_i$  (e.g., if  $g$  is  $\text{Parity}_k$ ), then the techniques we used for our general lower bound Theorem 1.3 may begin to apply.

However, we wish to close by drawing attention to a peculiar intermediate case. Suppose  $g$  is the following symmetric  $k$ -junta:

$$g(x) = \begin{cases} 0 & \text{if } \sum_{i=1}^k x_i = 1 \text{ or } 3 \leq \sum_{i=1}^k x_i < k/2, \\ 1 & \text{if } \sum_{i=1}^k x_i \in \{0, 2\} \text{ or } k/2 \leq \sum_{i=1}^k x_i \leq k, \end{cases}$$

In this case  $g$  is  $o_k(1)$ -close to  $\text{Maj}_k$ . Hence it would seem at first blush that the few jumps between 0 and 1 when  $\sum_{i=1}^k x_i \leq 3$  are irrelevant, and we should have an  $\omega_k(1)$  lower bound for testing isomorphism to this  $g$ .

But oddly, this is not clear. Because  $g$  is so close to  $\text{Maj}_k$ , it seems we would need to use ‘‘ $\mathcal{F}_{\text{no}}$  functions’’ which are fairly different from  $\text{Maj}_k$ , like the weighted majority function  $\text{WgtMaj}_k$  introduced in Section IV. However there is a *one-query* test that distinguishes between a function isomorphic to the above  $g$  and a function isomorphic to  $\text{WgtMaj}_k$ : simply query the string  $(0, 0, \dots, 0)$ , which has

<sup>2</sup>I.e., is the query complexity of the problem bounded below by  $\Omega(\log^* \ell)$ ? Or by  $\Omega(\log \ell)$ ? Or even by  $\Omega(\text{poly}(\ell))$ ?

value 1 under  $g$  and value 0 under  $\text{WgtMaj}_k$ ! We could fix this by changing  $\text{WgtMaj}_k(0, \dots, 0)$  to 0, but there are still problems. For example, the tester could query random strings having a  $1/k$  fraction of 1's. For such strings  $x$ ,  $\Pr[g(x) = 1]$  will be noticeably higher than  $\Pr[\text{Maj}_k(x) = 1]$ , because there is a good chance that the string  $x$  will contain exactly two 1's among the  $k$  coordinates on which  $g$  depends.

So strangely, even though strings in  $\{0, 1\}^k$  with zero, one, or two 1's constitute only an  $o_k(1) \ll \epsilon$  probability mass, a clever tester can exploit them for its advantage. This makes proving untestability for functions isomorphic to the above  $g$  somewhat tricky, and we leave it as a problem for future research.

#### ACKNOWLEDGMENTS

The second author would like to thank Adi Akavia and Guy Kindler for several helpful discussions. Both authors also thank the anonymous referees for valuable feedback on an earlier draft of this paper.

#### REFERENCES

- [1] Noga Alon, Eldar Fischer, Ilan Newman, and Asaf Shapira. A combinatorial characterization of the testable graph properties: It's all about regularity. In *Proc. of the 38th annual ACM Symposium on the Theory of Computing*, pages 251–260, 2006.
- [2] Noga Alon and Asaf Shapira. A characterization of the (natural) graph properties testable with one-sided error. In *Proc. of the 46th IEEE Symposium on Foundations of Computer Science*, pages 429–438, 2005.
- [3] Noga Alon and Asaf Shapira. Every monotone graph property is testable. In *Proc. of the 46th IEEE Symposium on Foundations of Computer Science*, pages 128–137, 2005.
- [4] Tim Austin and Terence Tao. On the testability and repair of hereditary hypergraph properties. To appear, *Random Struct. Alg.*
- [5] Mihir Bellare, Oded Goldreich, and Madhu Sudan. Free bits, PCPs and non-approximability – towards tight results. *SIAM J. Comput.*, 27(3):804–915, 1998.
- [6] Eldar Fischer. The difficulty of testing for isomorphism against a graph that is given in advance. *SIAM J. Comput.*, 34(5):1147–1158, 2005.
- [7] Eldar Fischer, Guy Kindler, Dana Ron, Shmuel Safra, and Alex Samorodnitsky. Testing juntas. *J. Comput. Syst. Sci.*, 68(4):753–787, 2004.
- [8] Eldar Fischer and Arie Matsliah. Testing graph isomorphism. *SIAM J. Comput.*, 38(1):207–225, 2008.
- [9] Parikshit Gopalan, Ryan O'Donnell, Yi Wu, and David Zuckerman. Fooling functions of halfspaces under product distributions. In *Proc. of the 25th annual IEEE Conference on Computational Complexity*, 2010.
- [10] Thomas Höglund. Sampling from a finite population. A remainder term estimate. *Scandinavian Journal of Statistics*, 5:69–71, 1978.
- [11] Kevin Matulef, Ryan O'Donnell, Ronitt Rubinfeld, and Rocco A. Servedio. Testing halfspaces. In *SODA '09: Proc. of the 19th Annual ACM -SIAM Symposium on Discrete Algorithms*, pages 256–264, 2009.
- [12] Kevin Matulef, Ryan O'Donnell, Ronitt Rubinfeld, and Rocco A. Servedio. Testing  $\pm 1$ -weight halfspace. In *RANDOM '09: Proc. of the 13th Intl. Workshop on Randomization and Computation*, pages 646–657, 2009.
- [13] Elchanan Mossel. Gaussian bounds for noise correlation of functions and tight analysis of long codes. In *Proc. of the 49th IEEE Symposium on Foundations of Computer Science*, pages 156–165, 2008.
- [14] Elchanan Mossel, Ryan O'Donnell, and Krzysztof Oleszkiewicz. Noise stability of functions with low influences: invariance and optimality. *Annals of Mathematics*, 171(1):295–341, 2010.
- [15] Michal Parnas, Dana Ron, and Alex Samorodnitsky. Testing basic boolean formulae. *SIAM J. Discrete Math.*, 16(1):20–46, 2002.
- [16] Yuval Peres. Noise stability of weighted majority, 2004. <http://arxiv.org/abs/math/0412377>.
- [17] Ronitt Rubinfeld and Madhu Sudan. Robust characterizations of polynomials with applications to program testing. *SIAM J. Comput.*, 25(2):252–271, 1996.
- [18] Andrew Chi-Chih Yao. Probabilistic computations: Toward a unified measure of complexity. In *Proc. of the 28th IEEE Symposium on Foundations of Computer Science*, pages 222–227, 1977.