

Communication Complexity of Statistical Distance

THOMAS WATSON, University of Memphis

We prove nearly matching upper and lower bounds on the randomized communication complexity of the following problem: Alice and Bob are each given a probability distribution over n elements, and they wish to estimate within $\pm\epsilon$ the statistical (total variation) distance between their distributions. For some range of parameters, there is up to a $\log n$ factor gap between the upper and lower bounds, and we identify a barrier to using information complexity techniques to improve the lower bound in this case. We also prove a side result that we discovered along the way: the randomized communication complexity of n -bit Majority composed with n -bit Greater Than is $\Theta(n \log n)$.

CCS Concepts: • Theory of computation → Communication complexity;

Additional Key Words and Phrases: Communication, complexity, statistical, distance

ACM Reference format:

Thomas Watson. 2018. Communication Complexity of Statistical Distance. *ACM Trans. Comput. Theory* 10, 1, Article 2 (January 2018), 11 pages.
<https://doi.org/10.1145/3170708>

1 INTRODUCTION

Statistical (a.k.a. total variation) distance is a standard measure of the distance between two probability distributions and is ubiquitous in theoretical computer science. Expressing the distributions (over a universe of n elements) as vectors of probabilities $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_n)$, the statistical distance is defined as

$$\Delta(x, y) := \frac{1}{2} \sum_{i \in [n]} |x_i - y_i| = \max_{S \subseteq [n]} \left| \sum_{i \in S} x_i - \sum_{i \in S} y_i \right| = \max_{S \subseteq [n]} \left(\sum_{i \in S} x_i - \sum_{i \in S} y_i \right).$$

This measure has various interpretations, such as the minimum over all couplings of the probability that the sample from x and the sample from y are unequal, or twice the maximum advantage an observer can achieve in guessing whether a random sample came from x or from y (where x or y is used with probability 1/2 each).

Given its pervasiveness, it is natural to inquire about the computational complexity of estimating the statistical distance between two distributions x and y that are given as input. This topic has been studied before in several contexts:

This work was supported by NSF grant CCF-1657377. A preliminary version was published as [32].

Author's address: T. Watson, Dunn Hall 315, 3725 Norriswood Avenue, University of Memphis, Memphis, TN, 38111; email: Thomas.Watson@memphis.edu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 ACM 1942-3454/2018/01-ART2 \$15.00

<https://doi.org/10.1145/3170708>

- [26] showed that when each of x and y is succinctly represented by an algorithm that takes uniform random bits and produces a sample from that distribution (so our actual input is the description of this pair of algorithms), then (a decision version of) the problem of estimating $\Delta(x, y)$ is complete for the complexity class SZK (statistical zero knowledge). (For results about the complexity of other problems where the inputs are succinctly represented distributions, see [4, 13–15, 30, 31].)
- [3, 10, 28] studied the complexity of statistical distance estimation when an algorithm is only given black-box access to oracles that produce samples from the distributions specified by x and y . (For results about the complexity of other problems where the inputs are black-box samples from distributions, see the surveys [8, 15, 25].)
- [11, 12] studied the space complexity of (a generalization of) statistical distance estimation when the vectors x and y are provided as data streams.

1.1 Communication Upper and Lower Bounds

We study the statistical distance estimation problem in the context of communication complexity: Alice is given the vector x , Bob is given the vector y , and they wish to output a value in the range $[\Delta(x, y) - \epsilon, \Delta(x, y) + \epsilon]$. We let $\text{STAT-DIST}_{n, \epsilon}$ denote this two-party search problem. For any two-party search problem F , we let $R(F)$ denote the minimum worst-case communication cost of any randomized protocol (allowing both public and private coins) such that for each input, the output is correct with probability at least $3/4$. (For our problem $\text{STAT-DIST}_{n, \epsilon}$, the $3/4$ can be replaced by any constant in the range $(1/2, 1)$ since we can amplify success probability by taking the median of multiple trials.) The following is a clean summary of our bounds.

$$\text{THEOREM 1.1. } R(\text{STAT-DIST}_{n, \epsilon}) \text{ is } \begin{cases} \Theta(1/\epsilon^2) & \text{if } 1 > \epsilon \geq 1/O(\sqrt{n}) \\ \Omega(n) \text{ and } O(n \log n) & \text{if } 1/\omega(\sqrt{n}) \geq \epsilon \geq 1/2^{o(n \log n)} \\ \Theta(\log(1/\epsilon)) & \text{if } 1/2^{\Omega(n \log n)} \geq \epsilon > 0 \end{cases}$$

We also go ahead and ascertain the deterministic communication complexity (denoted with D instead of R) of this problem. We prove Theorem 1.1 and Theorem 1.2 in Section 2.

THEOREM 1.2. $D(\text{STAT-DIST}_{n, \epsilon}) = \Theta(n \log(1/\epsilon))$ provided ϵ is at most a sufficiently small constant.

Closing the gap in Theorem 1.1 is a principal open problem. We get slightly better bounds in certain narrow ranges of ϵ (see the proof), but, e.g., it remains open to prove our conjecture that $R(\text{STAT-DIST}_{n, 1/2^n}) \geq \omega(n)$. A natural strategy is to use information complexity lower-bound techniques; however, in Section 3, we exhibit a barrier to accomplishing this. Specifically, for a large class of inputs having a certain type of product structure (which arises naturally from attempts to use the direct sum property of information complexity), and for a wide range of ϵ , $\text{STAT-DIST}_{n, \epsilon}$ can be solved with $O(n)$ information cost and 0 error probability. This suggests that to improve the $\Omega(n)$ bound, we may need to look at inputs not having the aforementioned product structure, and we are at a loss for techniques in this case.

1.2 Composing with Majority

We take this opportunity to prove other results that we discovered in the process of trying to analyze $\text{STAT-DIST}_{n, \epsilon}$. Recall the famous direct sum conjecture stating that computing k independent copies of a two-party function should require $\Omega(k)$ times as much randomized communication as computing one copy. A somewhat stronger version of the conjecture states that even just comput-

ing the AND of k independent copies should still require $\Omega(k)$ times as much communication.¹ [16] proved the query complexity analog of this AND-composition conjecture, as well as a communication complexity version that is weaker than the full conjecture in two senses: it is *qualitatively* weaker since instead of converting a protocol for AND_k composed with F into a plain randomized (BPP-type) protocol for F with factor $\Omega(k)$ savings, the conversion results in a protocol in a slightly stronger model (which has been variously called 2WAPP [16, 17], two-sided smooth rectangle bound [19], and relaxed partition bound [20]); it is *quantitatively* weaker since besides the $\Omega(k)$ savings, the conversion incurs a logarithmic additive loss due to the use of the “information odometer” of [6]. (We provide the precise statement in Section 4.)

We prove that when composing with the k -bit Majority function MAJ_k instead of AND_k , the above quantitative deficiency can be avoided: we get a perfect $\Omega(k)$ factor savings by circumventing the need for the odometer (although we retain the qualitative deficiency). For the applications in [1, 16], the logarithmic additive loss in the AND-composition result was immaterial albeit perhaps a slight nuisance. In some settings, however, that loss would be damaging; one such setting is the following corollary (which holds by combining our MAJ-composition result with the lower bound of [7] for the Greater-Than function GT_n on n -bit inputs).

THEOREM 1.3. $R(\text{MAJ}_n \circ \text{GT}_n^n) = \Theta(n \log n)$.

Evaluating the function $\text{MAJ}_n \circ \text{GT}_n^n$ can be described by a story: Alice and Bob have taken some exams and know their own scores, and they wish to determine the victor of their rivalry: who got a higher score on the most exams?

We prove the MAJ-composition result and provide details about Theorem 1.3 in Section 4. We make the stronger conjecture that Theorem 1.3 should hold even with AND_n instead of MAJ_n ; this would follow from an $\Omega(\log n)$ information complexity lower bound for GT_n with respect to a distribution only over 1-inputs (which is open but may be doable).

1.3 Preliminaries

We define AND_n , Or_n , and MAJ_n as the And, Or, and Majority functions on n bits, and EQ_n , GT_n , DISJ_n , and GH_n as the Equality, Greater-Than, Set-Disjointness, and Gap-Hamming two-party functions where Alice and Bob each get n bits. We use \mathbb{P} for probability, \mathbb{E} for expectation, \mathbb{H} for Shannon entropy, and \mathbb{I} for mutual information. We generally use uppercase letters for random variables and corresponding lowercase letters for particular outcomes.

Randomized protocols by default have both public and private coins. We let $CC(\Pi)$ denote the worst-case communication cost of protocol Π . We let $IC_D(\Pi) := \mathbb{I}(T ; X | Y, R) + \mathbb{I}(T ; Y | X, R)$ denote the (internal) information cost with respect to (X, Y) sampled from the input distribution D , where the random variables T and R represent the communication transcript and public coins of Π , respectively.

2 COMMUNICATION UPPER AND LOWER BOUNDS

We now prove Theorem 1.1 and Theorem 1.2. As a preliminary technicality, we note that for the upper bounds, we may assume that each of the probabilities x_i and y_i can be written exactly in binary with $\log(n/\epsilon) + O(1)$ bits. This is because if we truncate the binary representations to that many bits and reassign the lost probability to an arbitrary element in both x and y , this ensures that at most $\epsilon/4$ mass has been shifted within each distribution, so their statistical distance changes by

¹More precisely, the complexity of the composed function should be at least $\Omega(k)$ times (complexity of original function $- O(1)$). The $-O(1)$ is necessary since, e.g., computing the AND of k independent copies of the 2-bit AND function still only needs $O(1)$ communication.

at most $\epsilon/2$; then, to obtain an ϵ -estimation for the original x and y , we can run a protocol to get an $(\epsilon/2)$ -estimation for the new x and y .

PROOF OF THEOREM 1.1. In fact, we show that $R(\text{STAT-DIST}_{n,\epsilon})$ is always

- (i) $O(1/\epsilon^2)$
- (ii) $O(\max(n \log n, \log(1/\epsilon)))$
- (iii) $\Omega(\min(1/\epsilon^2, n))$
- (iv) $\Omega(\log(1/\epsilon))$,

which gives a slightly more detailed picture than the statement of Theorem 1.1.

The proof of (i) is inspired by the “correlated sampling lemma” that has been used in the context of parallel repetition [18, 23, 24] and earlier in the context of LP rounding [21]. As noted above, we may assume that each probability x_i and y_i is a multiple of $1/m$ for some integer $m := O(n/\epsilon)$. We make use of an $O(1)$ -communication equality testing protocol that accepts with probability 1 when the inputs are equal and accepts with probability exactly $1/2$ when the inputs are unequal (e.g., by using the inputs to index into a uniformly random public string and comparing the bits at those indices).

Here is the protocol witnessing (i). Alice and Bob repeat the following $O(1/\epsilon^2)$ times:

- Publicly sample a uniformly random ordering of $[n] \times [m]$.
- Alice finds the first (i_A, j_A) in the ordering such that $x_{i_A} \geq j_A/m$.
- Bob finds the first (i_B, j_B) in the ordering such that $y_{i_B} \geq j_B/m$.
- Run the equality test on (i_A, j_A) and (i_B, j_B) .

Then they output $q/(1-q)$, where

$$q := \min(1/2, \text{fraction of iterations where equality test rejected}).$$

To analyze the correctness, let $\delta := \Delta(x, y)$ and let p denote the probability the equality test rejects in a single iteration of the loop. We claim that $p = \delta/(1 + \delta)$ (and hence $\delta = p/(1 - p)$). To see this, define the following subsets of $[n] \times [m]$: $A := \{(i, j) : x_i \geq j/m \text{ and } y_i < j/m\}$, $B := \{(i, j) : x_i < j/m \text{ and } y_i \geq j/m\}$, and $C := \{(i, j) : x_i \geq j/m \text{ and } y_i \geq j/m\}$. Then $|A| = |B| = \delta m$ and $|C| = (1 - \delta)m$. The first (i^*, j^*) in the ordering to land in $A \cup B \cup C$ is uniformly distributed in that set. Thus, with probability $\delta/(1 + \delta)$, we have $(i^*, j^*) \in A$, in which case $(i_A, j_A) = (i^*, j^*) \neq (i_B, j_B)$; and with probability $\delta/(1 + \delta)$, we have $(i^*, j^*) \in B$, in which case $(i_A, j_A) \neq (i^*, j^*) = (i_B, j_B)$; and with probability $(1 - \delta)/(1 + \delta)$, we have $(i^*, j^*) \in C$, in which case $(i_A, j_A) = (i^*, j^*) = (i_B, j_B)$. It follows that the equality test rejects with probability $\frac{\delta}{1+\delta} \cdot \frac{1}{2} + \frac{\delta}{1+\delta} \cdot \frac{1}{2} + \frac{1-\delta}{1+\delta} \cdot 0 = \delta/(1 + \delta)$.

By a Chernoff bound, the number of iterations guarantees that with probability at least $3/4$, $|q - p| \leq \epsilon/8$. Since $\frac{d}{dp}[p/(1 - p)] = 1/(1 - p)^2 \in [1, 4]$ for all $p \in [0, 1/2]$, it follows that $|\text{output} - \delta| = |q/(1 - q) - p/(1 - p)| \leq \epsilon/2$ whenever $|q - p| \leq \epsilon/8$ and $q \in [0, 1/2]$. This proves (i).

To prove (ii), we exploit the fact that the Greater-Than function GT_k with k -bit inputs can be computed with error probability $\gamma > 0$ and $O(\log(k/\gamma))$ bits of communication (by running the standard binary-search-based protocol [22, p. 170] for $O(\log(k/\gamma))$ many steps). As noted above, we may assume each probability x_i and y_i has $\log(n/\epsilon) + O(1)$ bits.

Here is the protocol witnessing (ii). For each $i \in [n]$, Alice and Bob compute $\text{GT}(x_i, y_i)$ with error probability $1/(4n)$. Then Alice sends Bob the sum of x_i over all i for which the protocol for $\text{GT}(x_i, y_i)$ accepted, and Bob sends Alice the sum of y_i over the same i s. They output Alice’s sum minus Bob’s sum. By a union bound, with probability at least $3/4$, each of the GT tests returns the correct answer, in which case the final output is correct by definition. The communication cost is $O(n \log(n \log(n/\epsilon)) + \log(n/\epsilon)) \leq O(\max(n \log n, \log(1/\epsilon)))$.

To prove (iii), we use a reduction from the Gap-Hamming partial function $\text{GH}_{n,\epsilon}$, in which the goal is to determine whether the relative Hamming distance between Alice’s and Bob’s length- n bit strings is $> 1/2 + \epsilon$ or $< 1/2 - \epsilon$. It is known that $R(\text{GH}_{n,\epsilon}) \geq \Omega(\min(1/\epsilon^2, n))$ [9, 27, 29]. Here is

the reduction: Alice transforms $a \in \{0, 1\}^n$ into a distribution x over $[2n]$ by letting $x_{2i-a_i} = 1/n$ for each $i \in [n]$ (and letting all other entries of x be 0). Bob transforms b into y in the same way. Then $\Delta(x, y)$ equals the relative Hamming distance between a and b , so a protocol for $\text{STAT-DIST}_{2n, \epsilon}$ can distinguish the two cases (by whether the output is above or below $1/2$).

To prove (iv), consider any correct randomized protocol for $\text{STAT-DIST}_{n, \epsilon}$, and fix any set of $1/(3\epsilon)$ many pairs of distributions having statistical distances $0, 3\epsilon, 6\epsilon, 9\epsilon, \dots$. There must exist some outcome of the randomness of the protocol such that the induced deterministic protocol is correct on at least three-fourths of those inputs. But then the same transcript cannot occur for any two of these $1/(4\epsilon)$ inputs since the statistical distances are more than 2ϵ apart. Thus, at least $1/(4\epsilon)$ transcripts are necessary, so the communication cost must be at least $\log(1/\epsilon) - 2$. \square

PROOF OF THEOREM 1.2. For the upper bound, assuming each probability x_i and y_i is a multiple of $1/m$ for some integer $m := O(n/\epsilon)$, we employ the trivial protocol where Alice sends a specification of her distribution to Bob (who then responds with the $(\log(n/\epsilon) + O(1))$ -bit answer). We just need to count the number of such distributions: $\binom{m+n-1}{n-1} \leq \left(\frac{e \cdot (m+n-1)}{n-1}\right)^{n-1} \leq (O(1/\epsilon))^n$. Hence, only $O(n \log(1/\epsilon))$ bits are needed to specify a distribution.

The proof of the lower bound is basically a Gilbert–Varshamov argument for codes in the Manhattan metric. Specifically, we claim that there is a set of $2^{\Omega(n \log(1/\epsilon))}$ many distributions over $[n]$ that pairwise have statistical distance $> 2\epsilon$. Then for any distinct distributions x and x' from this set, the inputs (x, x) and (x', x') cannot share the same transcript in any correct protocol for $\text{STAT-DIST}_{n, \epsilon}$, because if they did then (x, x') would also share that transcript, but (x, x) requires output $\leq \epsilon$ while (x, x') requires output $> \epsilon$. Hence, any correct protocol has at least $2^{\Omega(n \log(1/\epsilon))}$ transcripts and so has communication cost $\Omega(n \log(1/\epsilon))$.

To see the claim, first note that the number of distributions whose probabilities are multiples of $1/m$ is $(\Omega(1/\epsilon))^n$, while the number of such distributions within statistical distance $\leq 2\epsilon$ of any fixed such distribution can be simply upper bounded by $2^n \cdot \binom{4\epsilon m + n}{n} \leq (O(1))^n$. Hence, if we keep greedily adding to a set any distribution that has statistical distance $> 2\epsilon$ from every distribution we picked so far, then the number of iterations this process can continue is at least $(\Omega(1/\epsilon))^n / (O(1))^n \geq (\Omega(1/\epsilon))^n$, which is $2^{\Omega(n \log(1/\epsilon))}$ provided ϵ is at most a sufficiently small constant. \square

3 INFORMATION UPPER BOUND

As motivation, we first note that $R(\text{STAT-DIST}_{3n, \epsilon}) \geq \Omega(n)$ for any $\epsilon < 1/(2n)$ follows by a reduction from the Set-Disjointness function DISJ_n (where the 1-inputs are pairs of length- n bit strings representing disjoint sets). Here is the reduction: Alice transforms $a \in \{0, 1\}^n$ into a distribution x over $[3n]$ by applying the following rule for each $i \in [n]$: if $a_i = 1$, then $x_{3i} = 1/n$, and if $a_i = 0$, then $x_{3i-1} = 1/n$. Bob uses the following rule to transform b into y : if $b_i = 1$, then $y_{3i} = 1/n$, and if $b_i = 0$, then $y_{3i-2} = 1/n$. (All other entries of x and y are set to 0.) Then $\Delta(x, y)$ equals the fraction of coordinates $i \in [n]$ such that $a_i = 0$ or $b_i = 0$, which is 1 if $\text{DISJ}_n(a, b) = 1$ and $\leq 1 - 1/n$ if $\text{DISJ}_n(a, b) = 0$. Thus, a protocol for $\text{STAT-DIST}_{3n, \epsilon}$ can distinguish the two cases (by whether the output is above or below $1 - 1/(2n)$).

The information complexity proof of the lower bound $R(\text{DISJ}_n) \geq \Omega(n)$ [2] shows that in a certain sense, the n coordinates (each of which is an AND_2 “gadget”) each contribute $\Omega(1)$ to the information cost, and these contributions add up over the coordinates. Thus, it is plausible that by similar reasoning, a lower bound of the form $R(\text{STAT-DIST}_{O(n), o(1/n)}) \geq \omega(n)$ could be shown by starting with an appropriate “gadget” that contributes $\omega(1)$ to the information cost. We now show that a very general formulation of this approach cannot work.

Let us examine more closely the instances (x, y) that arise from the above reduction from DISJ_n . The $3n$ coordinates are grouped into blocks of size 3, and within each block, Alice's and Bob's distributions both have probability exactly $1/n$, and conditioned on the block, they have statistical distance either 0 or 1 (so the block contributes either 0 or $1/n$ to the statistical distance of the whole input). This can be viewed as a “product structure” that enables the blocks to be considered independently of each other, and allows the contributions of the blocks to be summed to get a lower bound on the information cost of a STAT-DIST protocol.

Let us formalize a general class of inputs having the above product structure. Suppose C is an arbitrary constant, and the distributions have Cn coordinates that are grouped into blocks of size C . Assume Alice's and Bob's distributions satisfy the following promise: within each block, they both have probability exactly $1/n$, and conditioned on the block, the statistical distance is in either $[\ell - 2\epsilon, \ell]$ or $[u, u + 2\epsilon]$ for some $\ell < u$ (so the block's contribution to the statistical distance of the whole input is $1/n$ times that conditioned statistical distance). We use $\text{STAT-DIST}_{Cn, \epsilon}^{\ell, u}$ to denote the partial function with this promise on the input.

Note that if $u - \ell > 4\epsilon n$, then a protocol for $\text{STAT-DIST}_{Cn, \epsilon}^{\ell, u}$ could be used to determine the fraction of blocks for which the conditioned statistical distance falls in the lower range versus the upper range. This could be useful in an attempt to prove a $\omega(n)$ bound using information complexity techniques, e.g., via our “Majority-composition” result (Theorem 4.2). However, such an attempt would be futile:

PROPOSITION 1. *If C is a constant and $u - \ell \geq \epsilon$, then there is a protocol Π solving $\text{STAT-DIST}_{Cn, \epsilon}^{\ell, u}$ with 0 error probability and such that $IC_D(\Pi) \leq O(n)$ holds for every distribution D over inputs.*

In fact, the proposition holds even if we allow a different ℓ, u for each block. Also, note that the support of D is allowed to include inputs that do not satisfy the promise.

PROOF OF PROPOSITION 1. It suffices to prove this for $n = 1$, since by [5, Theorem 4.2] we can run such a protocol on each block to estimate the conditioned statistical distance within ϵ . The average (over the blocks) of those estimates will be within ϵ of the statistical distance of the whole input, and the information cost just adds up over the n blocks.

Assuming $n = 1$, it suffices to determine whether $\Delta(x, y)$ falls in the lower range (outputting $\ell - \epsilon$ if so) or the upper range (outputting $u + \epsilon$ if so). Let $\gamma := 1/2^{\lceil \log(1/\epsilon) \rceil} \in (\epsilon/2, \epsilon]$, and keep in mind the intervals $[0, \gamma]$, $[\gamma, 2\gamma]$, $[2\gamma, 3\gamma]$, \dots . We make use of the fact that there exists an equality testing protocol with 0 error probability that has $O(1)$ information cost under every distribution [5, Section 3.4].

Here is our protocol for $n = 1$. Alice and Bob repeat the following for each nonempty $S \subseteq [C]$:

- Alice finds the integer k_A such that $\sum_{i \in S} x_i \in [k_A \gamma, (k_A + 1)\gamma]$.
- Bob finds the integer k_B such that $\sum_{i \in S} y_i + u \in [k_B \gamma, (k_B + 1)\gamma]$.
- For each $m \in \{0, 1, 2, 3, 4\}$, run the equality test on k_A and $k_B + m$.

If any of the $(2^C - 1) \cdot 5$ equality tests accept, then they output $u + \epsilon$; otherwise, they output $\ell - \epsilon$.

We argue correctness. If $\Delta(x, y) \in [u, u + 2\epsilon]$, then there exists a nonempty S such that $\sum_{i \in S} x_i$ is contained in the range $\sum_{i \in S} y_i + [u, u + 2\epsilon]$, which is a subset of $\bigcup_{m \in \{0, 1, 2, 3, 4\}} [(k_B + m)\gamma, (k_B + m + 1)\gamma]$ (since $\gamma > \epsilon/2$); hence, $k_A = k_B + m$ for some $m \in \{0, 1, 2, 3, 4\}$ and so one of the equality tests accepts. If $\Delta(x, y) \in [\ell - 2\epsilon, \ell]$, then for every nonempty S , we have $\sum_{i \in S} x_i \leq \sum_{i \in S} y_i + \ell$, so $k_A < k_B$ must hold since otherwise $[k_A \gamma, (k_A + 1)\gamma]$ would contain both $\sum_{i \in S} y_i + \ell$ and $\sum_{i \in S} y_i + u$ (contradicting $u - \ell \geq \epsilon \geq \gamma$); hence, all the equality tests reject.

As for the information cost, fix an arbitrary distribution over inputs. Each of the equality tests has $O(1)$ information cost (using the simple fact that the information cost is unaffected by Alice

and Bob applying deterministic functions to their inputs to obtain the inputs to the equality test). Then again, by [5, Theorem 4.2], we can simply sum up this $O(1)$ information cost over the $O(1)$ many equality tests (noting that although [5, Theorem 4.2] is stated for tasks applied to separate inputs, arbitrary correlations are allowed between the inputs, so the upper bound still holds if we have multiple tasks applied to the same input). \square

4 COMPOSING WITH MAJORITY

In this section, we follow a convention that has become common in recent literature: for a two-party (possibly partial) function $F : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}$ and a complexity class name C , we let $C(F)$ denote the minimum worst-case cost of any protocol for F in the model corresponding to C , and we also use C to denote the class of (families of) F s such that $C(F) \leq \text{polylog}(n)$. In particular, $\text{BPP}(F)$ is an alias for the plain randomized communication complexity $\text{R}(F)$ in the case of $\{0, 1\}$ -valued F , but we use the complexity class notation now for aesthetic consistency. We also need the following “2-sided WAPP” model.²

Definition 1. $2\text{WAPP}_\epsilon(F) := \min(\text{CC}(\Pi) + \log(1/\alpha))$ over all $\alpha > 0$ and protocols Π with output values $\{0, 1, \perp\}$ such that for all (x, y) , $\mathbb{P}[\Pi(x, y) \neq \perp] \leq \alpha$ and $\mathbb{P}[\Pi(x, y) = F(x, y)] \geq (1 - \epsilon)\alpha$.

For all F and constants $0 < \epsilon < 1/2$, we have $O(\text{BPP}(F)) \geq 2\text{WAPP}_\epsilon(F) \geq \Omega(\text{PP}(F))$, and thus $\text{BPP} \subseteq 2\text{WAPP}_\epsilon \subseteq \text{PP}$. It is not necessary to recall the communication complexity definition of PP , but we remark that 2WAPP_ϵ feels intuitively much closer to BPP , since there are many interesting classes sandwiched between 2WAPP_ϵ and PP [17]. The following is due to [17].

THEOREM 4.1 (AND-COMPOSITION). *For all F, k , and constants $0 < \epsilon < 1/2$, we have*

$$2\text{WAPP}_\epsilon(F) \leq O(\text{BPP}(\text{AND}_k \circ F^k)/k + \log \text{BPP}(\text{AND}_k \circ F^k)).$$

We prove that by using MAJ_k instead of AND_k , the logarithmic term can be avoided.

THEOREM 4.2 (MAJ-COMPOSITION). *For all F, k , and constants $0 < \epsilon < 1/2$, we have*

$$2\text{WAPP}_\epsilon(F) \leq O(\text{BPP}(\text{MAJ}_k \circ F^k)/k + 1).$$

PROOF OF THEOREM 1.3. As noted in the proof of Theorem 1.1, GT_n has a protocol with error probability $1/(4n)$ and communication cost $O(\log n)$. By running this on each of n coordinates, with probability at least $3/4$ all the outputs will be correct, so a protocol witnessing $\text{BPP}(\text{MAJ}_n \circ \text{GT}_n^n) \leq O(n \log n)$ can be obtained by applying MAJ_n to all these outputs. The matching lower bound follows by combining Theorem 4.2 with the result that $\text{PP}(\text{GT}_n) \geq \Omega(\log n)$ [7]. \square

Theorem 4.2 follows by stringing together the following three lemmas. For any input distribution D (over the domain of F), we define the distributions $D^b := (D \mid F^{-1}(b))$ for $b \in \{0, 1\}$. We say a protocol Π is δ -correct for F if and only if $\mathbb{P}[\Pi(x, y) = F(x, y)] \geq 1 - \delta$ for all (x, y) .

LEMMA 4.3. *Fix any F, k , $0 < \delta < 1/2$, and input distribution D . For every δ -correct protocol Π for $\text{MAJ}_k \circ F^k$, there exists a δ -correct protocol Π' for F such that $\text{IC}_{D^b}(\Pi') \leq O(\text{CC}(\Pi)/k)$ holds for both $b \in \{0, 1\}$.*

LEMMA 4.4. *Fix any F , input distribution D , and protocol Π (not necessarily correct). Then*

$$\text{IC}_D(\Pi) - 4 \leq \sum_b \mathbb{P}_D[F^{-1}(b)] \cdot \text{IC}_{D^b}(\Pi) \leq \text{IC}_D(\Pi).$$

²There are two ways to define this model, which are equivalent up to a factor of 2 in ϵ . Our way was also used in [17] and is the same as the relaxed partition bound [20]. In [16], a “starred” notation was used for this, while the notation 2WAPP was reserved for the other definition, which is the same as the two-sided smooth rectangle bound [19].

LEMMA 4.5. Fix any F , constants $0 < \delta < \epsilon < 1/2$, and value c . If for every input distribution D there exists a δ -correct protocol Π for F such that $IC_D(\Pi) \leq c$, then $2\text{WAPP}_\epsilon(F) \leq O(c + 1)$.

Only the first inequality in Theorem 4.4 is needed for Theorem 4.2. Theorem 4.5 is due to [20]. Before we commence with the proofs of Theorem 4.3 and Theorem 4.4, we recall the following standard fact; see [5, Section 2.1] for a proof. (We apologize for overloading the D notation between this fact and the above lemmas, but there should be no confusion.)

FACT 1. Let A, B, C, D be four random variables. Then

- (i) $\mathbb{I}(A ; B | C) \leq \mathbb{I}(A ; B | C, D)$ if $\mathbb{I}(B ; D | C) = 0$;
- (ii) $\mathbb{I}(A ; B | C) \geq \mathbb{I}(A ; B | C, D)$ if $\mathbb{I}(B ; D | A, C) = 0$.

PROOF OF THEOREM 4.3. Assume k is odd for convenience. Consider a probability space with the following random variables: $Z \in \{0, 1\}^k$ is a uniformly random string of Hamming weight $\lceil k/2 \rceil$, $S := \{i : Z_i = 1\}$, (X, Y) is such that $(X_i, Y_i) \sim D^{Z_i}$ for each $i \in [k]$ independently, and T and R are the communication transcript and public coins, respectively, of Π on input (X, Y) . We use the subscript notation $X_{<i}$ and $X_{>i}$ for restrictions to coordinates in $\{1, \dots, i-1\}$ and $\{i+1, \dots, k\}$, we use the superscript notation X^S and X^{-S} for restrictions to coordinates in S and $[k] \setminus S$, and we may combine these so, e.g., $X_{>i}^{-S}$ is the restriction to coordinates in $\{i+1, \dots, k\} \setminus S$. We use corresponding notation for restrictions of Y . We have

$$\begin{aligned} & 2 \cdot CC(\Pi) \\ & \geq \mathbb{I}(T ; X^S | X^{-S}, Y, R, S) + \mathbb{I}(T ; Y^S | Y^{-S}, X, R, S) \\ & = \mathbb{E}_{s \sim S} \left[\sum_{i \in s} \mathbb{I}(T ; X_i | X_{<i}^s, X_{>i}^{-s}, Y, R, s) + \sum_{i \in s} \mathbb{I}(T ; Y_i | Y_{>i}^s, Y_{>i}^{-s}, X, R, s) \right] \\ & \geq \mathbb{E}_{s \sim S} \left[\sum_{i \in s} \mathbb{I}(T ; X_i | Y_i, X_{<i}, Y_{>i}, R, s) + \sum_{i \in s} \mathbb{I}(T ; Y_i | X_i, Y_{>i}, X_{<i}, R, s) \right] \\ & = \lceil k/2 \rceil \cdot \mathbb{E}_{\substack{s \sim S, \\ i \sim s, \\ r \sim R \\ x_{<i} \sim X_{<i}, \\ y_{>i} \sim Y_{>i}}} [\mathbb{I}(T ; X_i | Y_i, x_{<i}, y_{>i}, r, s) + \mathbb{I}(T ; Y_i | X_i, x_{<i}, y_{>i}, r, s)], \end{aligned}$$

where the second line is by the chain rule; the third line is by Fact 1.(i) since $X_{>i}^{-S}, Y_{<i}$ is independent of X_i given $Y_i, X_{<i}, Y_{>i}, R, s$ and since $Y_{<i}^{-S}, X_{>i}$ is independent of Y_i given $X_i, Y_{>i}, X_{<i}, R, s$; and $i \sim s$ on the fourth line means i is sampled uniformly at random from the set s .

Note that sampling $s \sim S$ and $i \sim s$ is equivalent to sampling $i \sim [k]$ and a uniformly random balanced bit string $z_{-i} \sim Z_{-i}$ indexed by $[k] \setminus \{i\}$ (and setting $z_i = 1$). We let $q \sim Q$ denote a sample of all the data $(i, z_{-i}, r, x_{<i}, y_{>i})$. In summary, we have

$$\mathbb{E}_{q \sim Q} [\mathbb{I}(T ; X_i | Y_i, q) + \mathbb{I}(T ; Y_i | X_i, q)] \leq (2/\lceil k/2 \rceil) \cdot CC(\Pi),$$

so by Markov's inequality, with probability $> 1/2$ over $q \sim Q$, we have

$$\mathbb{I}(T ; X_i | Y_i, q) + \mathbb{I}(T ; Y_i | X_i, q) \leq (4/\lceil k/2 \rceil) \cdot CC(\Pi), \quad (1)$$

where $(X_i, Y_i) \sim D^1$. By symmetric reasoning (interchanging the roles of 0 and 1), with probability $> 1/2$ over $q \sim Q$, Equation (1) also holds if we instead have $(X_i, Y_i) \sim D^0$. Thus, there exists a q (which we fix henceforth) such that Equation (1) holds both when $(X_i, Y_i) \sim D^1$ and when $(X_i, Y_i) \sim D^0$ (and in either case, $(X_j, Y_j) \sim D^{z_j}$ for $j \neq i$).

Now consider the protocol Π' where the input is interpreted as (x_i, y_i) , Alice privately samples $x_{>i} \sim (X_{>i} | y_{>i}, z_{>i})$, Bob privately samples $y_{<i} \sim (Y_{<i} | x_{<i}, z_{<i})$, and they run Π on the combined input (x, y) with public coins r . The conclusion of the previous paragraph is exactly

that $IC_{D^b}(\Pi') \leq (4/\lceil k/2 \rceil) \cdot CC(\Pi) \leq O(CC(\Pi)/k)$ holds for both $b \in \{0, 1\}$. Furthermore, Π' is δ -correct since Π is δ -correct and $F(x_i, y_i) = (\text{MAJ}_k \circ F^k)(x, y)$ with probability 1, for every (x_i, y_i) in F 's domain. \square

PROOF OF THEOREM 4.4. Consider a probability space with the following random variables: $(X, Y) \sim D$, $F := F(X, Y)$, and T and R are the communication transcript and public coins, respectively, of Π on input (X, Y) . Then we have

$$\begin{aligned} IC_D(\Pi) &= \mathbb{I}(T ; X | Y, R) + \mathbb{I}(T ; Y | X, R) \\ \sum_b \mathbb{P}_D[F^{-1}(b)] \cdot IC_{D^b}(\Pi) &= \mathbb{I}(T ; X | Y, R, F) + \mathbb{I}(T ; Y | X, R, F), \end{aligned}$$

and so the second inequality of Theorem 4.4 holds by Fact 1. (ii) since conditioned on X, Y, R , there is no remaining entropy in F and hence it is independent of T .

For the first inequality, we use the following result proven in [16].

LEMMA 4.6. *There exist numbers $c_{x,y}, c'_{x,y} \geq 0$ for each input (x, y) in the domain of F , such that*

- $IC_D(\Pi) = \mathbb{E}[c_{X,Y}]$,
- $IC_{D^b}(\Pi) = \mathbb{E}[c'_{X,Y} | F = b]$ for both $b \in \{0, 1\}$,
- for each (x, y) in the domain of F , letting $b := F(x, y)$, we have

$$c_{x,y} \leq c'_{x,y} + \log(1/\mathbb{P}[F = b | y]) + \log(1/\mathbb{P}[F = b | x]).$$

Hence, letting $p_{x,y} := \mathbb{P}[(X, Y) = (x, y)]$, we have

$$\begin{aligned} IC_D(\Pi) &= \sum_{(x,y)} p_{x,y} \cdot c_{x,y} \\ &\leq \sum_b \sum_{(x,y) \in F^{-1}(b)} p_{x,y} \cdot (c'_{x,y} + \log(1/\mathbb{P}[F = b | y]) + \log(1/\mathbb{P}[F = b | x])) \\ &= \sum_b \mathbb{P}[F = b] \cdot IC_{D^b}(\Pi) \\ &\quad + \sum_b \sum_{(x,y) \in F^{-1}(b)} p_{x,y} \cdot (\log(1/\mathbb{P}[F = b | y]) + \log(1/\mathbb{P}[F = b | x])). \end{aligned}$$

We claim that for both $b \in \{0, 1\}$, we have $\sum_{(x,y) \in F^{-1}(b)} p_{x,y} \cdot \log(1/\mathbb{P}[F = b | y]) \leq 1$ and $\sum_{(x,y) \in F^{-1}(b)} p_{x,y} \cdot \log(1/\mathbb{P}[F = b | x]) \leq 1$; it then follows that $IC_D(\Pi) \leq \sum_b \mathbb{P}[F = b] \cdot IC_{D^b}(\Pi) + 4$.

We just argue the claim for $b = 1$ and conditioning on y ; the other three cases are completely analogous. For $a \in \{0, 1\}$, define $p_y^a := \mathbb{P}[F = a \text{ and } Y = y] = \sum_{x : (x,y) \in F^{-1}(a)} p_{x,y}$. Then we have

$$\begin{aligned} \sum_{(x,y) \in F^{-1}(1)} p_{x,y} \cdot \log(1/\mathbb{P}[F = 1 | y]) &= \sum_y p_y^1 \cdot \log\left(\left(p_y^0 + p_y^1\right)/p_y^1\right) \\ &\leq \sum_y p_y^1 \cdot \left(\left(p_y^0 + p_y^1\right)/p_y^1\right) \\ &= 1. \end{aligned}$$

This finishes the proof. \square

ACKNOWLEDGMENTS

I thank Mika Göös for discussions and anonymous reviewers for comments.

REFERENCES

- [1] Anurag Anshu, Aleksandrs Belovs, Shalev Ben-David, Mika Göös, Rahul Jain, Robin Kothari, Troy Lee, and Miklos Santha. 2016. Separations in communication complexity using cheat sheets and information complexity. In *Proceedings of the 57th Symposium on Foundations of Computer Science (FOCS'16)*. IEEE, 555–564. DOI : <http://dx.doi.org/10.1109/FOCS.2016.66>
- [2] Ziv Bar-Yossef, T. S. Jayram, Ravi Kumar, and D. Sivakumar. 2004. An information statistics approach to data stream and communication complexity. *J. Comput. System Sci.* 68, 4 (2004), 702–732. DOI : <http://dx.doi.org/10.1016/j.jcss.2003.11.006>
- [3] Tugkan Batu, Lance Fortnow, Ronitt Rubinfeld, Warren Smith, and Patrick White. 2013. Testing closeness of discrete distributions. *J. ACM* 60, 1 (2013), 4. DOI : <http://dx.doi.org/10.1145/2432622.2432626>
- [4] Andrej Bogdanov, Elchanan Mossel, and Salil Vadhan. 2008. The complexity of distinguishing Markov random fields. In *Proceedings of the 12th International Workshop on Randomization and Computation (RANDOM'08)*. Springer, 331–342. DOI : http://dx.doi.org/10.1007/978-3-540-85363-3_27
- [5] Mark Braverman. 2015. Interactive information complexity. *SIAM J. Comput.* 44, 6 (2015), 1698–1739. DOI : <http://dx.doi.org/10.1137/130938517>
- [6] Mark Braverman and Omri Weinstein. 2015. An interactive information odometer and applications. In *Proceedings of the 47th Symposium on Theory of Computing (STOC'15)*. ACM, 341–350. DOI : <http://dx.doi.org/10.1145/2746539.2746548>
- [7] Mark Braverman and Omri Weinstein. 2016. A discrepancy lower bound for information complexity. *Algorithmica* 76, 3 (2016), 846–864. DOI : <http://dx.doi.org/10.1007/s00453-015-0093-8>
- [8] Clément Canonne. 2015. *A Survey on Distribution Testing: Your Data is Big. But Is It Blue?* Technical Report TR15-063. Electronic Colloquium on Computational Complexity (ECCC). Retrieved from <http://eccc.hpi-web.de/report/2015/063>.
- [9] Amit Chakrabarti and Oded Regev. 2012. An optimal lower bound on the communication complexity of Gap-Hamming-Distance. *SIAM J. Comput.* 41, 5 (2012), 1299–1317. DOI : <http://dx.doi.org/10.1137/120861072>
- [10] Siu On Chan, Ilias Diakonikolas, Paul Valiant, and Gregory Valiant. 2014. Optimal algorithms for testing closeness of discrete distributions. In *Proceedings of the 25th Symposium on Discrete Algorithms (SODA'14)*. ACM-SIAM, 1193–1203. DOI : <http://dx.doi.org/10.1137/1.9781611973402.88>
- [11] Joan Feigenbaum, Sampath Kannan, Martin Strauss, and Mahesh Viswanathan. 2002. An approximate L^1 -difference algorithm for massive data streams. *SIAM J. Comput.* 32, 1 (2002), 131–151. DOI : <http://dx.doi.org/10.1137/S0097539799361701>
- [12] Jessica Fong and Martin Strauss. 2001. An approximate L^p -difference algorithm for massive data streams. *Discrete Math. Theor. Comput. Sci.* 4, 2 (2001), 301–322.
- [13] Oded Goldreich, Amit Sahai, and Salil Vadhan. 1999. Can statistical zero knowledge be made non-interactive? or On the relationship of SZK and NISZK. In *Proceedings of the 19th International Cryptology Conference (CRYPTO'99)*. Springer, 467–484. DOI : http://dx.doi.org/10.1007/3-540-48405-1_30
- [14] Oded Goldreich and Salil Vadhan. 1999. Comparing entropies in statistical zero-knowledge with applications to the structure of SZK. In *Proceedings of the 14th Conference on Computational Complexity (CCC'99)*. IEEE, 54–73. DOI : <http://dx.doi.org/10.1109/CCC.1999.766262>
- [15] Oded Goldreich and Salil Vadhan. 2011. On the complexity of computational problems regarding distributions. *Studies in Complexity and Cryptography*. (2011), 390–405. DOI : http://dx.doi.org/10.1007/978-3-642-22670-0_27
- [16] Mika Göös, T. S. Jayram, Toniann Pitassi, and Thomas Watson. 2017. Randomized communication vs. partition number. In *Proceedings of the 44th International Colloquium on Automata, Languages, and Programming (ICALP)*. Schloss Dagstuhl, 52:1–52:15. DOI : <http://dx.doi.org/10.4230/LIPIcs.ICALP.2017.52>
- [17] Mika Göös, Shachar Lovett, Raghu Meka, Thomas Watson, and David Zuckerman. 2016. Rectangles are nonnegative juntas. *SIAM J. Comput.* 45, 5 (2016), 1835–1869. DOI : <http://dx.doi.org/10.1137/15M103145X>
- [18] Thomas Holenstein. 2009. Parallel repetition: Simplification and the no-signaling case. *Theory Comput.* 5, 1 (2009), 141–172. DOI : <http://dx.doi.org/10.4086/toc.2009.v005a008>
- [19] Rahul Jain and Hartmut Klauck. 2010. The partition bound for classical communication complexity and query complexity. In *Proceedings of the 25th Conference on Computational Complexity (CCC'10)*. IEEE, 247–258. DOI : <http://dx.doi.org/10.1109/CCC.2010.31>
- [20] Iordanis Kerenidis, Sophie Laplante, Virginie Lerays, Jérémie Roland, and David Xiao. 2015. Lower bounds on information complexity via zero-communication protocols and applications. *SIAM J. Comput.* 44, 5 (2015), 1550–1572. DOI : <http://dx.doi.org/10.1137/130928273>
- [21] Jon Kleinberg and Éva Tardos. 2002. Approximation algorithms for classification problems with pairwise relationships: Metric labeling and Markov random fields. *J. ACM* 49, 5 (2002), 616–639. DOI : <http://dx.doi.org/10.1145/585265.585268>

- [22] Eyal Kushilevitz and Noam Nisan. 1997. *Communication Complexity*. Cambridge University Press.
- [23] Anup Rao. 2011. Parallel repetition in projection games and a concentration bound. *SIAM J. Comput.* 40, 6 (2011), 1871–1891. DOI: <http://dx.doi.org/10.1137/080734042>
- [24] Ran Raz. 2011. A counterexample to strong parallel repetition. *SIAM J. Comput.* 40, 3 (2011), 771–777. DOI: <http://dx.doi.org/10.1137/090747270>
- [25] Ronitt Rubinfeld. 2012. Taming big probability distributions. *ACM Crossroads* 19, 1 (2012), 24–28. DOI: <http://dx.doi.org/10.1145/2331042.2331052>
- [26] Amit Sahai and Salil Vadhan. 2003. A complete problem for statistical zero knowledge. *J. ACM* 50, 2 (2003), 196–249. DOI: <http://dx.doi.org/10.1145/636865.636868>
- [27] Alexander Sherstov. 2012. The communication complexity of Gap Hamming Distance. *Theory Comput.* 8, 1 (2012), 197–208. DOI: <http://dx.doi.org/10.4086/toc.2012.v008a008>
- [28] Paul Valiant. 2011. Testing symmetric properties of distributions. *SIAM J. Comput.* 40, 6 (2011), 1927–1968. DOI: <http://dx.doi.org/10.1137/080734066>
- [29] Thomas Vidick. 2012. A concentration inequality for the overlap of a vector on a large set, with application to the communication complexity of the Gap-Hamming-Distance problem. *Chicago J. Theoret. Comput. Sci.* 2012, 1 (2012), 1–12. DOI: <http://dx.doi.org/10.4086/cjcs.2012.001>
- [30] Thomas Watson. 2015. The complexity of deciding statistical properties of samplable distributions. *Theory Comput.* 11 (2015), 1–34. DOI: <http://dx.doi.org/10.4086/toc.2015.v011a001>
- [31] Thomas Watson. 2016. The complexity of estimating min-entropy. *Comput. Complex.* 25, 1 (2016), 153–175. DOI: <http://dx.doi.org/10.1007/s00037-014-0091-2>
- [32] Thomas Watson. 2017. Communication complexity of statistical distance. In *Proceedings of the 21st International Workshop on Randomization and Computation (RANDOM'17)*. Schloss Dagstuhl, 49:1–49:10. DOI: <http://dx.doi.org/10.4230/LIPIcs.APPROX-RANDOM.2017.49>

Received June 2017; revised October 2017; accepted October 2017