

One-sided error communication complexity of Gap Hamming Distance

Egor Klenin ^{*1} and Alexander Kozachinskiy ^{†1, 2}

¹Lomonosov Moscow State University

²National Research University Higher School of Economics

Abstract

Assume that Alice has a binary string x and Bob a binary string y , both strings are of length n . Their goal is to output 0, if x and y are at least L -close in Hamming distance, and output 1, if x and y are at least U -far in Hamming distance, where $L < U$ are some integer parameters known to both parties. If the Hamming distance between x and y lies in the interval (L, U) , they are allowed to output anything. This problem is called the Gap Hamming Distance. In this paper we study public-coin one-sided error communication complexity of this problem. The error with probability at most $1/2$ is allowed only for pairs at Hamming distance at least U . In this paper we determine this complexity up to factors logarithmic in L . The protocol we construct for the upper bound is simultaneous.

1 Communication complexity of GHD

Given two strings $x = x_1 \dots x_n \in \{0, 1\}^n$, $y = y_1 \dots y_n \in \{0, 1\}^n$, Hamming distance between x and y is defined as the number of positions, where x and y differ:

$$d(x, y) = |\{i \in \{1, \dots, n\} \mid x_i \neq y_i\}|.$$

Let $L < U \leq n$ be integer numbers. In this paper we consider the following communication problem $\text{GHD}_{L,U}$, called the Gap Hamming Distance problem:

*yegorklenin@gmail.com

†kozlach@mail.ru

Definition 1. Let Alice receive an n -bit string x and Bob an n -bit string y such that either $d(x, y) \leq L$, or $d(x, y) \geq U$. They have to output 0, if the first inequality holds, and 1, if the second inequality holds. If the promise is not fulfilled, they may output anything.

1.1 Prior work

1.1.1 Two-sided error upper bounds

Let $R(\text{GHD}_{L,U})$ denote randomized two-sided error public coin communication complexity of $\text{GHD}_{L,U}$. It is known (see [11, 15]) that $R(\text{GHD}_{L,U}) = O(L^2/(U - L)^2)$ (assuming the constant error probability less than $1/2$).

The paper [9] established the upper bound $R(\text{GHD}_{L,U}) = O(L \log L)$ in the case $U = L + 1$, that is, there is no gap. This bound is much better than $O(L^2/(U - L)^2)$, which is $O(L^2)$ in this case.

It turns out that the protocols attaining these two upper bounds are *simultaneous*. That is, in these protocols Alice and Bob do not communicate at all, but rather send messages to the third party, Charlie, who then computes the output of the protocol. Charlie doesn't see inputs of Alice and Bob but sees public coins. The corresponding model, called simultaneous message passing (SMP) model, is even more restricted than one-way public-coin communication: every simultaneous protocol can be converted into one-way protocol without increasing communication cost.

1.1.2 One-sided error public coin communication protocols

The one-sided error public coin communication complexity will be denoted by R^0 . The superscript 0 means that the protocol is allowed to err only for input pairs which are at least U -far in Hamming distance. Here we assume that the maximal probability of error¹ is $1/2$. The superscript 1 will mean the opposite: protocols are allowed to err only for input pairs which are at least L -close in Hamming distance.

Let us first note that for all x, y we have

$$\text{GHD}_{L,U}(x, y) = \neg \text{GHD}_{n-U, n-L}(\neg x, y).$$

(Alice flips all bits of her input string.) Thus $\text{GHD}_{L,U}(x, y)$ reduces to $\text{GHD}_{n-U, n-L}$ and the other way around. This reduction maps 0-instances

¹by the standard amplification argument we could have any constant between 0 and 1 instead of $1/2$.

to 1-instances and vice versa. This observation implies that

$$R^1(\text{GHD}_{L,U}) = R^0(\text{GHD}_{n-U,n-L}).$$

Thus it suffices to study only one of these quantities and we will stick to R^0 (the error is allowed when the distance is at least U).

The paper [8] noticed that for all $U > L$ it holds that $R^0(\text{GHD}_{L,U}) = O(L \log n)$. Once again, there is a public coin SMP protocol attaining this bound (which is just a simple modification of the standard protocol for EQUALITY).

1.1.3 GHD and the lower bounds in data streams and property testing

Several works used GHD to obtain lower bounds for streaming algorithms and for property testing problems. As was discovered in [14] by Woodruff, there is a reduction from GHD to a number of fundamental data stream problems, including the problem of estimating frequency moments. More specifically, if there is a $\Omega(n)$ lower bound against any r -round two-sided error communication protocol for $\text{GHD}_{n/2-\Theta(\sqrt{n}), n/2+\Theta(\sqrt{n})}$, then there is a $\Omega(1/(r\varepsilon^2))$ lower bound on the space complexity of any r -pass streaming algorithm estimating the frequency moments in a data stream within a factor of $(1 + \varepsilon)$.

In [14] Woodruff proved $\Omega(n)$ lower bound against 1-round protocols (see also [10] for more direct and simple proof). In a subsequent works ([2, 3]) this lower bound was extended to $O(1)$ -round protocols. Finally, $\Omega(n)$ lower bound in the most general setting, when there is no restriction on the number of rounds at all, was obtained in [6, 13, 12].

As it turns out, lower bounds on the one-sided error version of GHD are also useful. In [5] Buhrman, Cleve and Wigderson proved that for any constant $c < 1$ it holds that $R^1(\text{GHD}_{0,cn}) = \Omega(n)$. Moreover, they showed that $\Omega(n)$ lower bound holds also for a weaker version of $\text{GHD}_{0,cn}$ problem, in which Hamming distance between the inputs is either 0 or *exactly* cn (provided that cn is an even integer).

Blais, Brody and Matulef in [1] used this result to obtain lower bounds on testing decision trees and signed majorities with one-sided error.

Further, Brody and Woodruff ([4]) used lower bound on one-sided error GHD from [5] to obtain lower bounds for streaming algorithms with *one-sided approximation*, i.e., for algorithms which either always return an overestimate or always return an underestimate on the objective function. Their results include lower bounds for the problem of over(under)-estimating the

number of non-zero rows in a matrix and the Earth Mover Distance between two multisets.

1.2 This work

In this paper we study public-coin one-sided error communication complexity R^0 of $\text{GHD}_{L,U}$. Once again, the error is allowed only for pairs at Hamming distance at least U .

1.2.1 The upper bound

Our main result is a one-sided error public-coin simultaneous protocol for $\text{GHD}_{L,U}$ on n -bit strings with communication complexity $O((L^2/U) \log L)$. It is constructed in the following 4 steps (description of the protocol in this section is a bit informal, and the precise bounds can be found below in the paper). Let us stress that steps 1 and 2 are enough to obtain $O((L^2/U) \log n)$ solution; the purpose of steps 3 and 4 is to replace $O(\log n)$ -factor by $O(\log L)$ -factor. Importance of eliminating dependency on n in the upper bounds was also acknowledged in previous works ([15, 9]).

Step 1. On this step we construct our main novel protocol, called *the Triangle Inequality Protocol*. This protocol communicates $O((L^2/U) \log n)$ bits (which is a bit more than required, since $\log L$ is replaced by $\log n$) and solves the $\text{GHD}_{L,U}$ problem when the ratio U/L is larger than a certain constant.

The protocol works as follows. It randomly splits x and y in $b = O(L^2/U)$ blocks x^1, \dots, x^b and y^1, \dots, y^b . The i th bit x_i of x goes in the block x^j where j is chosen at random with uniform probability distribution over $\{1, \dots, b\}$, and decisions for different i 's are independent. Each bit y_i of y goes in the block y^j with the same index as x_i goes in. This partition is made using the shared random source (so that the parties have the same partition). Both parties also read random strings r^1, \dots, r^b from the shared random source and Alice communicates $d(x^j, r^j)$ to Charlie for all $j = 1, \dots, b$. Bob does the same with $d(y^j, r^j)$. Thus the communication is $b \log n = O((L^2/U) \log n)$. Charlie outputs 0 if the sum

$$\sum_{j=1}^b |d(x^j, r^j) - d(y^j, r^j)|$$

is at most L and 1 otherwise. By the triangle inequality each term in this sum is at most $d(x^j, y^j)$ and thus the sum is at most $d(x, y)$. Therefore this protocol does not err if $d(x, y) \leq L$.

On the other hand, if $d(x, y) \geq U \geq C'L$ for a certain constant C' , then for any fixed j the average value of $d(x^j, y^j)$ is at least 2. From the properties of binomial distributions it follows that we have $d(x^j, y^j) \geq d(x, y)/10b$ with probability at least $1/3$. The value $d(x^j, r^j) - d(y^j, r^j)$ is distributed as the distance from the origin in a random walk with $d(x^j, y^j)$ steps along a line (each step has length 1 and is directed to the left or to the right with equal probabilities). From the properties of random walks it follows that for every j we have $|d(x^j, r^j) - d(y^j, r^j)| > \sqrt{d(x^j, y^j)}$ with constant positive probability. These two facts imply that with constant probability the sum $\sum_{j=1}^b |d(x^j, r^j) - d(y^j, r^j)|$ is $\Omega(b\sqrt{d(x, y)/10b}) = \Omega(\sqrt{bd(x, y)})$. Recall that $b = O(L^2/U)$ and we assume that $d(x, y) \geq U$. If the constant hidden in O -notation is large enough then the lower bound $\Omega(\sqrt{bd(x, y)})$ for the sum $\sum_{j=1}^b |d(x^j, r^j) - d(y^j, r^j)|$ is larger than L .

Step 2. In [8] it was noticed that for all $L < U$ there is one-sided error public-coin simultaneous protocol for $\text{GHD}_{L,U}$ with communication $O(L \log n)$. This protocol is just a modification of the standard protocol for equality and it never errs for inputs at distance at most L .

Our second protocol runs the Triangle Inequality Protocol if $U > C'L$ and the protocol from [8] otherwise. Notice that in the latter case $L = O(L^2/U)$, and thus we obtain a protocol with communication $O((L^2/U) \log n)$ for all L, U .

Step 3. On this step we use the techniques from [9] to replace the $\log n$ factor by a $\log L$ factor. More specifically, we run the protocol from the second step for the strings u, v of length $O(L^8)$ obtained from the original strings x, y by the following transformation. As in the Triangle Inequality Protocol, we split x, y into $b = O(L^8)$ blocks and then replace each block by the parity of its bits. Obviously, $d(u, v) \leq d(x, y)$. We then show that $d(u, v) = d(x, y)$ with constant probability provided $d(x, y) \leq L^4$. Therefore this protocol has constant one-sided error probability for all input pairs with $d(x, y) \leq L^4$. By construction this protocol communicates $O((L^2/U) \log L)$ bits.

Step 4. Finally, to handle the case $d(x, y) > L^4$, we consider the following protocol. We run the protocol from step 3 and then a simplified version of the Triangle Inequality Protocol. If any of these two protocols output 1, we output 1 and otherwise 0. The simplified version of the Triangle Inequality Protocol works as follows. Alice and Bob read a random n -bit string r from the shared random source. They compute distance from their inputs to r . Observe that due to triangle inequality $|d(x, r) - d(y, r)| \leq d(x, y)$. Hence $d(x, y) \leq L$ implies $|d(x, r) - d(y, r)| \leq L$. On the other hand, if

$d(x, y) > L^4$, then due to the properties of random walks with constant positive probability it holds that $|d(x, r) - d(y, r)| > L^2$.

Thus step 4 is reduced to the following communication problem. Alice holds a number $a \in \{0, 1, \dots, n\}$, Bob holds a number $b \in \{0, 1, \dots, n\}$ and it is known that either $|a - b| \leq L$ or $|a - b| > L^2$. The goal is to find out whether the first or the second inequality is true. We construct a public-coin simultaneous protocol with communication $O(\log L)$ which always outputs 0 when $|a - b| \leq L$ and which with some constant positive probability outputs 1 when $|a - b| > L^2$.

There is a simple SMP protocol communicating $O(\log L + \log \log n)$ bits to solve even a gap-less (L vs $L + 1$) version of this problem. Let p_1, \dots, p_k be the first $k = (4L + 2) \cdot \log_2(n + L)$ primes. Parties read a random $p \in \{p_1, \dots, p_k\}$ from the shared random source. Alice and Bob communicate $a \pmod p$ and $b \pmod p$ to Charlie. He checks, whether there is $i \in [-L, L]$ such that $a + i \pmod p = b \pmod p$. If there is such i , he outputs 0, otherwise 1. A straightforward analysis shows that this protocol has an error at most $1/2$ only in the case when $|a - b| \geq L + 1$.

The problem with this protocol is that a and b range from 0 to n . To get rid of $O(\log \log n)$ term we do the following. Instead of taking remainders modulo p_1, \dots, p_k we just hash our $O(n)$ -size universe into $O(L)$ -size universe simply by taking remainder modulo $4L + 2$. Of course this may lead to a collision when two number which were far from each other become L -close. We resolve this issue by considering $Z_0 + \dots + Z_a$ and $Z_0 + \dots + Z_b$ instead of a and b , where Z_0, \dots, Z_n are independent symmetric Bernoulli random variables. It can be shown that provided $|a - b| > L^2$, the difference $Z_{a+1} + \dots + Z_b$ is distributed almost uniformly modulo $4L + 2$. This guaranties that the collision probability is by a constant bounded away from 1.

1.2.2 Lower bounds

As it turns out, a very simple argument proves an almost matching $\Omega(L^2/U)$ lower bound. We include this argument for completeness.

As we mentioned, provided that U is even and $U = (1 - \Omega(1))n$, the paper [5] proves $\Omega(n)$ lower bound on one-sided error communication complexity R^1 of an easier version of $\text{GHD}_{0,U}$, in which the distance between inputs is either 0 or exactly U . However, we need a lower bound in the regime when U is very close to n . We observe that a simple modification of a proof from [5] works as well in such regime when one switches to a harder problem, in which the distance between inputs can be greater than U . Namely, we show that $R^1(\text{GHD}_{0,U}) = \Omega((n - U)^2/n)$ for $\text{GHD}_{0,U}$ on n -bit strings.

As a corollary we obtain the lower bound $\Omega(L^2/U)$ for one-sided error complexity R^0 of $\text{GHD}_{L,U}$ (the error is allowed when the distance is at least U). As we explained earlier, R^1 of $\text{GHD}_{0,U-L}$ on U -bit strings equals R^0 of $\text{GHD}_{L,U}$ on U -bit strings. As the former is $\Omega((U - (U - L))^2/U)$ we obtain the lower bound $\Omega(L^2/U)$ for the latter. On the other hand, the problem $\text{GHD}_{L,U}$ on U -bit strings reduces to the problem $\text{GHD}_{L,U}$ on n -bit strings (Alice and Bob append $n - U$ zeros to their strings), hence the one-sided complexity R^0 of the latter is also $\Omega(L^2/U)$.

1.2.3 The summary

Let us summarize our results.

Theorem 1. *The one-sided error public-coin communication complexity R^0 of $\text{GHD}_{L,U}$ on n -bit strings is at most*

$$O\left(\left(\frac{L^2}{U} + 1\right) \log(L + 2)\right)$$

(The error is allowed only when the distance is at least U .) There is a public-coin simultaneous protocol attaining this bound.

Theorem 2. *The one-sided error public-coin communication complexity R^1 of $\text{GHD}_{0,U}$ on n -bit strings is at least*

$$\Omega\left(\frac{(n - U)^2}{n} + 1\right).$$

(The error is allowed only when the distance is 0.)

Corollary 3. *The one-sided error public-coin communication complexity R^0 of $\text{GHD}_{L,U}$ on n -bit strings is at least*

$$\Omega\left(\frac{L^2}{U} + 1\right).$$

(The error is allowed only when the distance is at least U .)

Thus our results determine the one-sided public-coin communication complexity of $\text{GHD}_{L,U}$ (up to a factor $O(\log L)$) in the case when the parties are allowed to err only for input pairs at distance at least U . If the parties are allowed to err only for input pairs at distance at most L , then the one-sided public-coin communication complexity of $\text{GHD}_{L,U}$ is $(n - U)^2/(n - L)$ up to a factor of $O(\log(n - U))$.

2 Preliminaries

2.1 Communication Complexity

Let $f : \mathcal{X} \times \mathcal{Y} \rightarrow \{0, 1\}$ be a Boolean function.

Definition 2. *A deterministic communication protocol is a rooted binary tree, in which each non-leaf vertex is associated either with Alice or with Bob and each leaf is labeled by 0 or 1. Each non-leaf vertex v , associated with Alice, is assigned a function $f_v : \mathcal{X} \rightarrow \{0, 1\}$ and each non-leaf vertex u , associated with Bob, is assigned a function $g_u : \mathcal{Y} \rightarrow \{0, 1\}$. For each non-leaf vertex one of its out-going edges is labeled by 0 and other one is labeled by 1.*

A computation according to a deterministic protocol runs as follows. Alice is given $x \in \mathcal{X}$, Bob is given $y \in \mathcal{Y}$. They start at the root of tree. If they are in a non-leaf vertex v , associated with Alice, Alice sends $f_v(x)$ to Bob and they move to the son of v by the edge labeled by $f_v(x)$. If they are in a non-leaf vertex, associated with Bob, they act in a similar same way, however this time it is Bob who sends a bit to Alice. When they reach a leaf, they output the bit which labels this leaf.

Definition 3. *Communication complexity of a deterministic protocol π , denoted by $CC(\pi)$, is defined as the depth of the corresponding binary tree.*

Randomized protocols with shared randomness (aka public-coin protocols) can be defined as follows:

Definition 4. *A public-coin communication protocol is a probability distribution over deterministic protocols. Communication complexity of a public-coin protocol τ , denoted by $CC(\tau)$, is defined as $\max_{\pi} CC(\pi)$, where π is taken over the deterministic protocols from the support of τ (recall that τ is a distribution).*

Given a public-coin protocol τ , Alice and Bob choose the deterministic protocol to be executed according to the distribution, defined by τ .

Definition 5. *We say that a public-coin protocol computes a partial function f with error probability ε , if for every pair of inputs (x, y) in the domain of f with probability at least $1 - \varepsilon$ that protocol outputs $f(x, y)$. Randomized communication complexity of f is defined as*

$$R_{\varepsilon}(f) = \min_{\pi} CC(\pi),$$

where minimum is over all protocols that compute f with error probability ε .

A deterministic *simultaneous* protocol τ is a triple $\langle \phi, \psi, \theta \rangle$ where

$$\begin{aligned}\phi : \mathcal{X} &\rightarrow \{0, 1\}^{c_1}, & \psi : \mathcal{Y} &\rightarrow \{0, 1\}^{c_2}, \\ \theta : \{0, 1\}^{c_1} \times \{0, 1\}^{c_2} &\rightarrow \{0, 1\}.\end{aligned}$$

The communication cost of τ is $c_1 + c_2$. A public-coin simultaneous protocol π is a probability distribution over deterministic simultaneous protocols. Communication cost of π is the maximal possible communication cost of τ , where τ is a deterministic simultaneous protocol taken from the support of π .

Assume that Alice is given $x \in \mathcal{X}$ and Bob is given $y \in \mathcal{Y}$. The output of a public-coin simultaneous protocol π on (x, y) is a random variable defined as follows. Sample a deterministic simultaneous protocol $\tau = \langle \phi, \psi, \theta \rangle$ according to π . Output $\theta(\phi(x), \psi(y))$.

If for $i \in \{0, 1\}$ we require that the protocol never errs on inputs from $f^{-1}(i)$, then the corresponding notion is called “randomized one-sided error communication complexity” and is denoted by $R_\varepsilon^i(f)$.

The Gap Hamming Distance problem is the problem of computing the following partial function:

$$\text{GHD}_{L,U}^n(x, y) = \begin{cases} 0 & d(x, y) \leq L, \\ 1 & d(x, y) \geq U, \\ \text{undefined} & L < d(x, y) < U, \end{cases} \quad \text{for } x, y \in \{0, 1\}^n.$$

2.2 Hamming Space

Definition 6. *The function*

$$h(x) = x \log_2 \frac{1}{x} + (1 - x) \log_2 \frac{1}{1 - x}$$

is called the Shannon function.

For any $B \subset \{0, 1\}^n$ define $\text{diam}(B) = \max_{x, y \in B} d(x, y)$. Let $V_2(n, r)$ denote the size of Hamming ball of radius r , that is $V_2(n, r) = \binom{n}{0} + \dots + \binom{n}{r}$.

We will use the following well-known facts about the size of Hamming balls.

Proposition 4 ([7]). *If $r \leq \frac{n}{2}$, then $V_2(n, r) \leq 2^{h(\frac{r}{n})n}$.*

Proposition 5 ([7]). *If $B \subset \{0, 1\}^n$, r is natural, $\text{diam}(B) \leq 2r$ and $n \geq 2r + 1$, then*

$$|B| \leq V_2(n, r).$$

Propositions 5, 4 and the fact that $h'(1/2) = 0, h''(1/2) < 0$ easily imply the following

Lemma 6. *Assume that $r < n/2$. Then the cardinality of every set $B \subset \{0, 1\}^n$ with $\text{diam}(B) \leq 2r$ is at most $2^{n(1-c(1-(2r/n))^2)}$ for some absolute positive constant c .*

2.3 Probability Theory

Definition 7 (Probability distributions). *Let $B(n, p)$ denote the binomial distribution with parameters $n \in \mathbb{N}$ and $p \in (0, 1)$. For every natural n let S_n denote the one-dimensional random walk with n steps. More specifically, let S_n be equal to*

$$S_n = X_1 + \dots + X_n,$$

where X_1, \dots, X_n are independent random variables taking values in $\{-1, 1\}$, such that for each i the following holds: $\Pr[X_i = 1] = \Pr[X_i = -1] = \frac{1}{2}$.

3 The upper bound

The protocol for Theorem 1 is a combination of three different protocols. The most important of them solves $\text{GHD}_{L,U}$ with one sided error in the case when U/L exceeds some constant. Its communication length is $O((L^2/U + 1) \log n)$. We call that protocol the ‘‘Triangle Inequality Protocol’’, because it uses the triangle inequality for Hamming distance.

3.1 The Triangle Inequality Protocol

The following Lemma is the standard fact of Probability Theory:

Lemma 7. *There exists a positive constant $\alpha > 0$ such that for every m it holds that*

$$\Pr[S_m \geq \sqrt{m}] \geq \alpha,$$

where S_m denotes one-dimensional random walk with m steps, i.e., S_m is equal to the sum of m independent random variables, each taking the values 1 and -1 with probabilities $1/2$.

Everywhere below α stands for the constant from Lemma 7.

Let $x, y \in \{0, 1\}^n$ denote Alice's and Bob's input strings, respectively. The parties set $b = \lceil CL^2/U + 1 \rceil$, where $C = 360/\alpha^2$. Then they use public coins to sample a function $\chi : \{1, 2, \dots, n\} \rightarrow \{1, 2, \dots, b\}$ uniformly at random. They use χ to divide x and y into b blocks

$$x^1, \dots, x^b, \quad y^1, \dots, y^b.$$

The block x^j consists of all bits x_i of x such that $\chi(i) = j$. Similarly, y^j consists of all bits y_i with $\chi(i) = j$. The order in which bits of j th block are arranged is not important, the parties care only that they use the same order.

Then they use public coins to sample b random strings r^1, \dots, r^b of the same lengths, as x^1, \dots, x^b and y^1, \dots, y^b . Alice then sends b numbers to Charlie:

$$d(x^1, r^1), \dots, d(x^b, r^b).$$

In turn, Bob sends

$$d(y^1, r^1), \dots, d(y^b, r^b).$$

Then Charlie computes the sum

$$T = \sum_{j=1}^b \left| d(x^j, r^j) - d(y^j, r^j) \right|.$$

If $T \leq L$, Charlie outputs 0. Otherwise he outputs 1.

If $d(x, y) \leq L$, then the protocol always outputs 0. Indeed, since Hamming distance satisfies the triangle inequality, we have that

$$T = \sum_{j=1}^b \left| d(x^j, r^j) - d(y^j, r^j) \right| \leq \sum_{j=1}^b d(x^j, y^j) = d(x, y) \leq L.$$

Thus this protocol has a one-sided error: it can err only if $d(x, y) \geq U$. Now we will estimate the probability of error in the case when $d(x, y) \geq U$.

Lemma 8. *Assume that $U \geq 2b$. Then the protocol for the input pair x, y at distance at least U outputs 1 with some positive constant probability (more specifically, with probability at least $\alpha/6$).*

Proof. Assume that U, L, x, y satisfy the assumption of the lemma. Fix $j = 1, \dots, b$. First we have to understand what is the distribution of the

random variable $|d(x^j, r^j) - d(y^j, r^j)|$. By construction Alice and Bob choose a random function χ that governs the partition of x, y into blocks. For each i such that $x_i \neq y_i$ the probability that x_i, y_i land into the block with number j is $1/b$. Hence the random variable $d(x^j, y^j)$ has binomial distribution $B(d(x, y), 1/b)$ with parameters $d(x, y)$ and $1/b$, i.e., the probability of the event $d(x^j, y^j) = k$ equals

$$\binom{d(x, y)}{k} (1/b)^k (1 - 1/b)^{d(x, y) - k}.$$

The average value of $d(x^j, y^j)$ is thus equal to $d(x, y)/b$.

Once x^j, y^j are determined, Alice and Bob sample r^j . The value $d(x^j, r^j) - d(y^j, r^j)$ can be represented as the sum of $|x^j| = |y^j|$ terms where each term corresponds to a number i with $\chi(i) = j$. If $x_i = y_i$ then the term is 0. Otherwise it is either -1 or 1 depending on whether the respective bit of r^j is equal to x_i or to y_i . Thus for every fixed partition into blocks the value $|d(x^j, r^j) - d(y^j, r^j)|$ is distributed as the distance from origin in the random walk along the line with $d(x^j, y^j)$ independent steps where each step is 1 with probability $1/2$ and -1 with the same probability.

To finish the proof we will use the following fact about binomial distribution.

Lemma 9. *If X is distributed according to the binomial distribution $B(n, p)$ and $pn \geq 2$, then*

$$\Pr \left[X > \frac{pn}{10} \right] \geq \frac{1}{3}.$$

Proof of Lemma 9. The expectation and variation of X are given by:

$$EX = pn, \quad \text{Var}X = p(1 - p)n \leq pn.$$

Hence by Chebyshev inequality we get

$$\Pr \left[X \leq \frac{pn}{10} \right] \leq \frac{\text{Var}X}{\left(pn \cdot \left(1 - \frac{1}{10} \right) \right)^2} \leq \frac{\frac{100}{81}}{pn} \leq \frac{100}{162} \leq \frac{2}{3}.$$

□

Recall that the random variable $d(x^j, y^j)$ has binomial distribution $B(d(x, y), 1/b)$ and we assume that $d(x, y)/b \geq U/b \geq 2$. Hence by Lemma 9 with probability at least $1/3$ we have $d(x^j, y^j) \geq d(x, y)/10b$.

Fix any partition into blocks such that $d(x^j, y^j) \geq d(x, y)/10b$. By Lemma 7 with probability at least α we have

$$|d(x^j, r^j) - d(y^j, r^j)| \geq \sqrt{d(x^j, y^j)} \geq \sqrt{d(x, y)/10b}.$$

We have proved that for every fixed j with probability at least $\alpha/3$ we have $|d(x^j, r^j) - d(y^j, r^j)| \geq \sqrt{d(x, y)/10b}$. A simple averaging argument shows that with probability at least $\alpha/6$ the fraction of j that satisfy this inequality is bigger than $\alpha/6$. Indeed, let the random variable θ denote the fraction of j that satisfy this inequality. Its average is at least $\alpha/3$. On the other hand, we can upperbound its average by the sum

$$\Pr[\theta > \alpha/6] \cdot 1 + \Pr[\theta \leq \alpha/6] \cdot (\alpha/6) \leq \Pr[\theta > \alpha/6] + \alpha/6.$$

Thus with probability $\alpha/6$ we have

$$\sum_{j=1}^b |d(x^j, r^j) - d(y^j, r^j)| > (\alpha/6)b\sqrt{d(x, y)/10b} = (\alpha/6)\sqrt{b \cdot d(x, y)/10}.$$

Recall that $b = \lceil CL^2/U + 1 \rceil$, where $C = 360/\alpha^2$, and $d(x, y) \geq U$. So the right hand side of the last displayed inequality is strictly larger than L . \square

If the ratio U/L is larger than a certain constant then the protocol solves $\text{GHD}_{L,U}$ with constant one-sided error-probability. One can verify that the assumption $U \geq 2b$ of Lemma 8 is met for all $U \geq 2CL + 4$.

Recall that the communication length of the protocol is $O((L^2/U + 1)\log n)$. Now we need a protocol with the same communication length for L, U such that $U \leq 2CL + 3$. Notice that in this case the upper bound $O((L^2/U)\log n)$ for communication boils down to $O(L\log n)$. A protocol with such performance was constructed in [8].

3.2 The protocol of [8]

For the reader's convenience and to stress that the protocol from [8] has one-sided error we give here its full description.

Here \oplus stands for the bit-wise XOR over n -bit vectors and $\langle \cdot, \cdot \rangle : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}$ denotes the inner product over \mathbb{F}_2 :

$$\langle a, b \rangle = \sum_{i=1}^n a_i b_i \pmod{2}.$$

Let $x, y \in \{0, 1\}^n$ denote Alice's and Bob's input strings, respectively. They use public coins to sample N vectors

$$R_1, \dots, R_N \in \{0, 1\}^n$$

independently uniformly at random. Alice sends $\langle x, R_1 \rangle, \dots, \langle x, R_N \rangle$ to Charlie. Bob does the same with y . If there is $f \in \{0, 1\}^n$ of Hamming weight at most L such that:

$$\langle x \oplus f, R_1 \rangle = \langle y, R_1 \rangle, \dots, \langle x \oplus f, R_N \rangle = \langle y, R_N \rangle, \quad (1)$$

then Charlie outputs 0. Otherwise Charlie outputs 1.

Such protocol costs $O(N)$ bits. If $d(x, y) \leq L$, then the protocol outputs 0 with probability 1. Indeed, $f = x \oplus y$ (which is of Hamming weight at most L in this case) satisfies (1).

Now assume that $d(x, y) > L$. Then any $f \in \{0, 1\}^n$ of Hamming weight at most L satisfies (1) only with probability at most 2^{-N} (because $x + f \neq y$). Hence the error probability of the protocol is at most $V_2(n, L) \cdot 2^{-N}$ in this case. Here $V_2(n, L)$ is the size of Hamming ball of radius L . As $V_2(n, L) \leq (n+1)^L$, it is enough to take $N = O(L \log n)$.

3.3 The simplified version of the Triangle Inequality Protocol

Thus for all L, U we have a public-coin simultaneous protocol with communication length $O((L^2/U + 1) \log n)$ to solve $\text{GHD}_{L,U}$ with constant one-sided error probability. To replace $\log n$ factor by $\log L$ factor we will need the following public-coin simultaneous protocol with communication length $O(\log L)$ to solve $\text{GHD}_{L, (4L+2+N_0)^4}$ with constant one-sided error probability. Here N_0 is a constant from the following Lemma.

Lemma 10. *There is a positive integer N_0 and a positive real c such that the following holds. Assume that m and N are positive integers and $N \geq \max\{N_0, m^2\}$. Consider N independent random variables Z_1, \dots, Z_N , where each variable takes the values 0 and 1 with probabilities 1/2. Then for every $i \in \{0, 1, \dots, m-1\}$ it holds that:*

$$\Pr[Z_1 + \dots + Z_N = i \pmod{m}] \geq \frac{c}{m}.$$

(the proof of this Lemma will be given in the end of this subsection). Notice that $O((L^2/U+1) \log L)$ becomes just $O(\log L)$ for $U = (4L+2+N_0)^4$.

The protocol. The parties use public coins to sample a vector $r \in \{0, 1\}^n$ uniformly at random. Alice and Bob compute the distance from r to their input strings. If $d(x, y) \leq L$, then by Triangle Inequality we have $|d(x, r) - d(y, r)| \leq d(x, y) \leq L$. On the other hand, assume that $d(x, y) \geq (4L + 2 + N_0)^4$. From Lemma 7 it follows that in this case with constant positive probability we have $|d(x, r) - d(y, r)| \geq \sqrt{d(x, y)} > (4L + 2)^2 + N_0^2$.

Consider the following auxiliary problem. Alice holds a number $a \in \{0, 1, \dots, n\}$, Bob holds a number $b \in \{0, 1, \dots, n\}$ and it is promised that either $|a - b| \leq L$ or $|a - b| > (4L + 2)^2 + N_0^2$. They want to know whether the first or the second inequality is true. As the previous paragraph shows, if there is a public-coin SMP protocol with communication length $O(\log L)$, which always outputs 0 when $|a - b| \leq L$ and which with constant positive probability outputs 1 when $|a - b| > (4L + 2)^2 + N_0^2$, then we are done.

Define $m = 4L + 2$. Use public coins to sample $n + 1$ independent random variables

$$Z_0, Z_1, Z_2, \dots, Z_n,$$

where each variable takes the values 0 and 1 with probabilities 1/2.

Alice sends $\sum_{i=0}^a Z_i \pmod{m}$ to Charlie, Bob sends $\sum_{i=0}^b Z_i \pmod{m}$ to Charlie. This takes only $O(\log m) = O(\log L)$ bits. Let (s, t) be any pair of integers satisfying the following three conditions:

$$s \equiv \sum_{i=0}^a Z_i \pmod{m} \tag{2}$$

$$t \equiv \sum_{i=0}^b Z_i \pmod{m} \tag{3}$$

$$|s - t| = \min \left\{ |s' - t'| : s' \equiv \sum_{i=0}^a Z_i \pmod{m}, t' \equiv \sum_{i=0}^b Z_i \pmod{m} \right\}. \tag{4}$$

Obviously, knowing $\sum_{i=0}^a Z_i \pmod{m}$, $\sum_{i=0}^b Z_i \pmod{m}$, Charlie is able to find (s, t) satisfying these three conditions. He then simply checks whether $|s - t| \leq L$. If this is the case, he outputs 0. Otherwise he outputs 1.

Once again, the protocol communicates only $O(\log L)$ bits, as required. Further, it is easy to see that the protocol has one-sided error. Indeed, assume that $|a - b| \leq L$. Note that a pair $(\sum_{i=0}^a Z_i, \sum_{i=0}^b Z_i)$ satisfies (2) and (3)

. Hence $|s - t| \leq \left| \sum_{i=0}^a Z_i - \sum_{i=0}^b Z_i \right| \leq |a - b| \leq L$.

Now, let's consider the case when $|a - b| > (4L + 2)^2 + N_0^2$. Assume without loss of generality that $a < b$. Let E be the event that there is no $r \in [-L, L]$ such that $Z_{a+1} + \dots + Z_b \equiv r \pmod{m}$. Let us verify that E implies that $|s - t| > L$ (which means that Charlie outputs 1). Indeed, observe that

$$t - s \equiv \sum_{i=0}^b Z_i - \sum_{i=0}^a Z_i \equiv Z_{a+1} + \dots + Z_b \pmod{m},$$

but if $|s - t| \leq L$, this contradicts E .

It only remains to show that E happens with constant positive probability. This follows from Lemma 10. Namely, this lemma implies that $\Pr[E] \geq \frac{c(m-2L-1)}{m} = c/2$. Parameters are chosen in such a way that restrictions of Lemma 10 are satisfied:

$$b - a > (4L + 2)^2 + N_0^2 \geq (\max\{N_0, 4L + 2\})^2 \geq \max\{N_0, m^2\}.$$

Proof of Lemma 10. Take N_0 to be the first natural satisfying the following condition: there exists $d > 0$ such that for all $N \geq N_0$ and for every k between $N/2 - \sqrt{N}$ and $N/2 + \sqrt{N}$ the following holds:

$$\Pr[Z_1 + \dots + Z_N = k] = \binom{N}{k} 2^{-N} \geq \frac{d}{\sqrt{N}}.$$

The existence of such N_0, d is just a standard corollary of the Stirling formula, applied to $\binom{N}{k}$.

Now let us show that for all $m > 0$, $N \geq m^2$ and $i \in \{0, 1, \dots, m-1\}$ the number of k between $N/2 - \sqrt{N}$ and $N/2 + \sqrt{N}$ such that $k \equiv i \pmod{m}$ is at least $\frac{\sqrt{N}}{m}$. The number of such k is equal to the number of $r \in \mathbb{Z}$ satisfying:

$$N/2 - \sqrt{N} \leq mr + i \leq N/2 + \sqrt{N},$$

This number is at least

$$\left\lfloor \frac{N/2 + \sqrt{N} - i}{m} \right\rfloor - \left\lfloor \frac{N/2 - \sqrt{N} - i}{m} \right\rfloor + 1 \geq \frac{2\sqrt{N}}{m} - 1.$$

Provided $N \geq m^2$, the last expression is at least $\frac{\sqrt{N}}{m}$.

Set $c = d$ and observe that for all m, N such that $m > 0$ and $N \geq \max\{N_0, m^2\}$ and for every $i \in \{0, 1, \dots, m-1\}$ it holds that

$$\Pr[Z_1 + \dots + Z_N \equiv i \pmod{m}] \geq \frac{\sqrt{N}}{m} \cdot \frac{d}{\sqrt{N}} = \frac{c}{m}.$$

□

3.4 The final protocol for Theorem 1

The protocol. Step 1. Alice and Bob first run the Simplified Triangle Inequality Protocol from the previous subsection. If that protocol outputs 1 they output 1 and halt. Otherwise they proceed to Step 2.

Step 2. They divide x and y into $w = 2(4L + 2 + N_0)^8$ blocks randomly (as in the construction of the Triangle Inequality Protocol). Let

$$x^1, \dots, x^w, \quad y^1, \dots, y^w$$

denote the resulting blocks. Let u_i be the XOR of all bits from x^i and let v_i be the XOR of all bits from y^i . Alice privately computes u_1, \dots, u_w and sets $u = u_1 \dots u_w$. Bob privately computes v_1, \dots, v_w and sets $v = v_1 \dots v_w$.

Recall that we have a protocol (a combination of the Triangle Inequality Protocol and the protocol of [8]) with communication length $O((L^2/U + 1) \log w) = O((L^2/U + 1) \log L)$ to solve $\text{GHD}_{L,U}$ on w -bit strings with constant positive one-sided error probability.

Alice and Bob run this protocol for input pair (u, v) (and not (x, y)). They output the result of this run.

The communication length of the constructed protocol is $O((L^2/U + 1) \log L)$. We have to show that it has one-sided constant error probability.

If $d(x, y) \leq L$ then the run of the Simplified Triangle Inequality Protocol will output 0 with probability 1. Thus they proceed to Step 2. The distance between u and v does not exceed the distance between x and y and hence is at most L . Thus the run of the second protocol also outputs 0 with probability 1.

Assume that $d(x, y) \geq U$. If $d(x, y) \geq (4L + 2 + N_0)^4$, then the Simplified Triangle Inequality Protocol outputs 1 with positive constant probability, they output 1 and halt.

Assume that $U \leq d(x, y) < (4L + 2 + N_0)^4$. We claim that in this case with constant positive probability we have $d(u, v) = d(x, y)$. Indeed, consider any two positions in which x and y differ. Those positions land

into the same block with probability $\frac{1}{w}$. By union bound, with probability at least

$$1 - \frac{d(x, y)^2}{w} \geq 1 - \frac{(4L + 2 + N_0)^8}{2(4L + 2 + N_0)^8} = 0.5$$

all the positions in which x and y differ land in different blocks. The latter means that for all i the blocks x^i and y^i differ in at most 1 position and hence $d(u, v) = d(x, y)$. Thus with probability at least $1/2$ we have $d(u, v) \geq U$ and Alice and Bob output 1 with positive constant probability on the second step.

4 The lower bound

In this section we prove Theorem 2.

Proof of Theorem 2. Let τ be a protocol witnessing $R_{\frac{1}{2}}^1(\text{GHD}_{0,U})$. Then the following hold:

- for each $x \in \{0, 1\}^n$ the protocol τ for input (x, x) outputs 0 with probability at least $\frac{1}{2}$;
- for all $x, y \in \{0, 1\}^n$ with $d(x, y) \geq U$ the protocol τ always outputs 1.

By the standard averaging argument due to von Neumann there is a deterministic protocol π such that

- the communication complexity of π is at most $R_{\frac{1}{2}}^1(\text{GHD}_{0,U})$;
- π outputs 0 for at least half of diagonal input pairs (x, x) ;
- π outputs 1 for all inputs pairs at Hamming distance at least U .

Consider any 0-leaf of π and the corresponding rectangle $R = A \times B \subset \{0, 1\}^n \times \{0, 1\}^n$. The number of diagonal pairs from R is equal to $|A \cap B|$. Diameter of $A \cap B$ must be less than U . Indeed, if there are $x, y \in A \cap B$ such that $d(x, y) \geq U$, then π outputs 0 for input pair (x, y) .

It turns out that the largest set of diameter $2r < n$ is the Hamming ball of radius r and the diameter of the latter is at most $2^{n(1-c(1-2r/n)^2)}$ for some positive constant c (Lemma 6).

Let $r = \lfloor U/2 \rfloor$. For $U = n$ the lower bound in Theorem 2 is constant and thus the statement is obvious. Therefore we may assume that $U < n$

and hence $r < n/2$. The diameter of $A \cap B$ is at most $2r$ (recall that the diameter of $A \cap B$ is strictly less than U). By Lemma 6 we have

$$|A \cap B| \leq 2^{n(1-c(1-2r/n)^2)} \leq 2^{n(1-c(1-U/n)^2)}.$$

We have shown that if R is the rectangle corresponding to a 0-leaf of π , then R covers at most $2^{n(1-c(1-U/n)^2)}$ diagonal pairs. As the total number of diagonal pairs covered by 0-leaves of π is at least 2^{n-1} , the number of 0-leaves in π is at least $2^{cn(1-U/n)^2-1}$. Thus we have

$$R_{\frac{1}{2}}^1(\text{GHD}_{0,U}) \geq c \cdot \frac{(n-U)^2}{n} - 1. \quad (5)$$

Obviously we also have

$$R_{\frac{1}{2}}^1(\text{GHD}_{0,U}) \geq 1. \quad (6)$$

From inequalities (5) and (6) we can easily deduce that

$$R_{\frac{1}{2}}^1(\text{GHD}_{0,U}) \geq \Omega\left(\frac{(n-U)^2}{n} + 1\right)$$

(for example, we can add these inequalities with appropriate positive weights). \square

References

- [1] BLAIS, E., BRODY, J., AND MATULEF, K. Property testing lower bounds via communication complexity. *Computational Complexity* 21, 2 (2012), 311–358.
- [2] BRODY, J., AND CHAKRABARTI, A. A multi-round communication lower bound for gap hamming and some consequences. In *Computational Complexity, 2009. CCC'09. 24th Annual IEEE Conference on* (2009), IEEE, pp. 358–368.
- [3] BRODY, J., CHAKRABARTI, A., REGEV, O., VIDICK, T., AND DE WOLF, R. Better gap-hamming lower bounds via better round elimination. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques*. Springer, 2010, pp. 476–489.
- [4] BRODY, J., AND WOODRUFF, D. P. Streaming algorithms with one-sided estimation. In *APPROX-RANDOM* (2011), Springer, pp. 436–447.

- [5] BUHRMAN, H., CLEVE, R., AND WIGDERSON, A. Quantum vs. classical communication and computation. In *Proceedings of the thirtieth annual ACM symposium on Theory of computing* (1998), ACM, pp. 63–68.
- [6] CHAKRABARTI, A., AND REGEV, O. An optimal lower bound on the communication complexity of gap-hamming-distance. *SIAM Journal on Computing* 41, 5 (2012), 1299–1317.
- [7] COHEN, G., HONKALA, I., LITSYN, S., AND LOBSTEIN, A. *Covering codes*, vol. 54. Elsevier, 1997.
- [8] GAVINSKY, D., KEMPE, J., AND DE WOLF, R. Quantum communication cannot simulate a public coin. *arXiv preprint quant-ph/0411051* (2004).
- [9] HUANG, W., SHI, Y., ZHANG, S., AND ZHU, Y. The communication complexity of the hamming distance problem. *Information Processing Letters* 99, 4 (2006), 149–153.
- [10] JAYRAM, T. S., KUMAR, R., AND SIVAKUMAR, D. The one-way communication complexity of hamming distance. *Theory of Computing* 4, 1 (2008), 129–135.
- [11] KUSHILEVITZ, E., OSTROVSKY, R., AND RABANI, Y. Efficient search for approximate nearest neighbor in high dimensional spaces. *SIAM Journal on Computing* 30, 2 (2000), 457–474.
- [12] SHERSTOV, A. A. The communication complexity of gap hamming distance. *Theory of Computing* 8, 1 (2012), 197–208.
- [13] VIDICK, T. 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 Journal of Theoretical Computer Science* 1 (2012).
- [14] WOODRUFF, D. Optimal space lower bounds for all frequency moments. In *Proceedings of the fifteenth annual ACM-SIAM symposium on Discrete algorithms* (2004), Society for Industrial and Applied Mathematics, pp. 167–175.
- [15] YAO, A. C.-C. On the power of quantum fingerprinting. In *Proceedings of the thirty-fifth annual ACM symposium on Theory of computing* (2003), ACM, pp. 77–81.