

Lower Bounds for Linear Locally Decodable Codes and Private Information Retrieval

Oded Goldreich*

Howard Karloff†

Leonard J. Schulman‡

Luca Trevisan§

Abstract

We prove that if a linear error-correcting code $C : \{0, 1\}^n \rightarrow \{0, 1\}^m$ is such that a bit of the message can be probabilistically reconstructed by looking at two entries of a corrupted codeword, then $m = 2^{\Omega(n)}$. We also present several extensions of this result.

We show a reduction from the complexity of one-round, information-theoretic Private Information Retrieval Systems (with two servers) to Locally Decodable Codes, and conclude that if all the servers' answers are linear combinations of the database content, then $t = \Omega(n/2^a)$, where t is the length of the user's query and a is the length of the servers' answers. Actually, 2^a can be replaced by $O(a^k)$, where k is the number of bit locations in the answer that are actually inspected in the reconstruction.

1 Introduction

This paper is concerned with two related notions. The first notion is that of locally decodable codes (LDC), which are error-correcting codes that allow recovery of individual information bits based on a few (randomly selected) codeword bits. The second notion is that of private information retrieval (PIR) schemes, which are protocols allowing users to retrieve desired data items from several (non-colluding) servers without yielding any information to any individual server. The relation between these notions has been observed by some researchers before, and is further established in this paper.

The study of LDCs was initiated by Katz and Trevisan [6], who established super-linear (but at most quadratic) lower bounds on the length of codes that allow recovery based on

*Computer Science Department, Weizmann Institute of Science, Rehovot, Israel. Supported by MINERVA Foundation, Germany. E-mail: oded@wisdom.weizmann.ac.il

†AT&T Labs—Research, USA. E-mail: howard@research.att.com

‡Caltech, USA. Supported in part by NSF grant 0049092 (previously 9876172), and by the Charles Lee Powell Foundation. E-mail: schulman@caltech.edu

§Computer Science Division, University of California, Berkeley. E-mail: luca@eecs.berkeley.edu

a constant number of bits. In contrast, the best known constructions of LDCs (supporting such efficient recovery) have sub-exponential length. This leaves a huge gap between the known lower and upper bounds, and an important research goal is to try to close this gap. We take a first step in this direction by closing the gap (via improved lower bounds) for the special case of *linear* LDCs in which recovery is based on *two* bits.

The study of PIR schemes was initiated by Chor, Goldreich, Kushilevitz and Sudan [4], who presented (among other schemes) a one-round, 2-server PIR scheme of communication complexity $O(n^{1/3})$. The question of whether their (2-server) PIR scheme has the lowest communication complexity possible has been open since. We present several results that are related to this question, where all our results relate to the special case of *one-round*, 2-server PIR schemes in which the servers' answers are always *linear* combinations of the data bits.

1.1 Locally Decodable Codes

In this paper we consider error-correcting codes with the following *local decodability* property: given a corrupted codeword it is possible to recover each bit of the original message by applying a probabilistic procedure that looks at only *two* entries of the corrupted codeword. The procedure should predict each bit with a constant advantage even when there is a constant fraction of errors in corrupted codeword. The Hadamard code satisfies this requirement, but unfortunately its codewords are exponentially longer than the message they encode. In this paper, we prove that this is essentially the best possible with respect to linear codes.

Let us first define formally the notion of a locally decodable code. For a natural number n , we let $[n] \stackrel{\text{def}}{=} \{1, \dots, n\}$. For $x \in \Sigma^m$ and $i \in [m]$, we let x_i be the i th element of x ; that is, $x = x_1 \cdots x_m$. For $y, z \in \Sigma^m$, we denote by $d(y, z)$ the number of locations on which y and z differ, that is, $d(y, z) = |\{i : y_i \neq z_i\}|$.

Definition 1.1 For reals δ, ϵ and an integer q , we say that $C : \Sigma^n \rightarrow \Gamma^m$ is a (q, δ, ϵ) -locally decodable code if there exists a probabilistic oracle machine A such that:

- In every invocation, A makes at most q queries (possibly adaptively). Query $i \in [m]$ to the oracle $y \in \Gamma^m$ is answered by y_i .
- For every $x \in \Sigma^n$, for every $y \in \Gamma^m$ with $d(y, \mathbf{C}(x)) \leq \delta m$, and for every $i \in [m]$, we have

$$\Pr[A^y(i) = x_i] \geq \frac{1}{2} + \epsilon,$$

where the probability is taken over the internal coin tosses of A .

An algorithm A satisfying the above requirements is called an (adaptive) (q, δ, ϵ) -local decoding algorithm for \mathbf{C} .

While it appears natural to allow adaptive reconstruction algorithms in our definition, we only know how to directly prove lower bounds in the non-adaptive case. Lower bounds for the non-adaptive case can be generalized to the adaptive case by using the following reduction.

Lemma 1.2 ([6]) *Let $\mathbf{C} : \Sigma^n \rightarrow \Gamma^m$ be an error-correcting code that has an adaptive $(2, \delta, \epsilon)$ -local decoding algorithm. Then \mathbf{C} also has a non-adaptive $(2, \delta, \epsilon/|\Gamma|)$ -local decoding algorithm.*

All the results that we will state (from now on) refer to non-adaptive reconstruction procedures, and “local decoding algorithm” and “locally decodable code” will always refer to the non-adaptive case. We omit the statement of the results for the adaptive case (which can be obtained by the application of the above lemma).

As stated above, our work focuses on *linear* codes. In particular, we will consider the following settings:

- $\Sigma = \Gamma = F$ is a finite field, and the function $\mathbf{C} : F^n \rightarrow F^m$ is a linear mapping between the vector spaces F^n and F^m . In Theorem 1.3 we deal with the special case $\Sigma = \Gamma = GF(2)$, while in Theorem 1.4 we deal with general fields.
- $\Sigma = \{0, 1\}$, $\Gamma = \{0, 1\}^l$, and $\mathbf{C} : \{0, 1\}^n \rightarrow \{0, 1\}^{lm}$ is linear. We deal with this case in Theorem 1.5.
- $\Sigma = \Gamma = \{0, 1\}^l$, and $\mathbf{C} : \{0, 1\}^{ln} \rightarrow \{0, 1\}^{lm}$ is linear. That is, we consider codes mapping a sequence of n blocks, each being a string of length l , to a sequence of m such blocks, and algorithms that recover a desired (entire) block by making two block-queries. We refer to such codes as block-block codes, and deal with them in Theorem 1.6.

Our main result is

Theorem 1.3 *Let $\Sigma = \Gamma = \{0, 1\}$, and let $\mathbf{C} : \Sigma^n \rightarrow \Gamma^m$ be a $(2, \delta, \epsilon)$ -locally decodable linear code. Then $m \geq 2^{\epsilon \delta n / 8}$.*

The result has the following extensions to larger alphabets (corresponding to the three cases discussed above). First, we consider an extension to linear codes over arbitrary finite fields.

Theorem 1.4 *Let $\mathbf{C} : F^n \rightarrow F^m$ be a $(2, \delta, \epsilon)$ -locally decodable linear code. Then $m \geq 2^{\frac{\epsilon \delta}{16} \cdot n - 1 - \log_2 |F|}$.*

Theorem 1.5 *Let $\mathbf{C} : \{0, 1\}^n \rightarrow (\{0, 1\}^l)^m$ be a $(2, \delta, \epsilon)$ -locally decodable linear code, and suppose that the decoder uses only k predetermined bits out of the l bits that it receives as answer to each query. Then $m \geq (1/f(k, l)) \cdot 2^{\epsilon \delta n / (8f(k, l))}$, where $f(k, l) = \sum_{i=0}^k \binom{l}{i} \leq \min\{2^l, 2l^k\}$.*

Theorem 1.6 *Let $\mathbf{C} : (\{0, 1\}^l)^n \rightarrow (\{0, 1\}^l)^m$ be a $(2, \delta, \epsilon)$ -locally decodable code that is a linear block-block code. Then $m \geq 2^{\frac{\epsilon \delta}{16} \cdot n - (\ell+1)^2}$.*

Theorem 1.4 is proved by an extension of the argument used in the proof of Theorem 1.3. Theorem 1.5 is proved by means of a reduction to the case $l = k = 1$ and an application of Theorem 1.3. Theorem 1.6 is proved by an extension of the argument used in the proof of Theorem 1.3.

1.2 Private Informational Retrieval

Loosely speaking, a Private Information Retrieval (PIR) scheme for k servers is a protocol by which a user can obtain the value of a desired bit out of n bits held by the servers without yielding the identity of this bit to any individual server (assuming that the servers do not cooperate in order to learn the identity of the desired bit). The aim is to obtain PIR schemes of low communication complexity (i.e., substantially lower than the obvious solution of having a server send all n bits to the user). We focus on one-round PIR schemes that are protocols in which the user sends a single message to each server, which responds also with a single message. In the definition below, Q represents the algorithm employed by the user to generate its queries, S_j represents the algorithm employed by the j th server, and R represents the recovery algorithm used by the user (once it gets the servers’ answers).

Definition 1.7 *A one-round, $(1 - \delta)$ -secure, 2-server PIR scheme for database size n , with recovery probability p , query size t and answer size a is a quadruple of deterministic algorithms $\bar{A} = (Q, S_1, S_2, R)$ with the following properties.*

Algorithmic operation: *On input $i \in [n]$ and (random-tape) $r \in \{0, 1\}^L$, algorithm Q outputs a pair of t -bit long queries; that is, $(q_1, q_2) \stackrel{\text{def}}{=} Q(i, r)$.*

On input a database $x \in \{0, 1\}^n$, and query $q \in \{0, 1\}^t$, algorithm S_1 (resp., S_2) returns an answer $S_1(x, q) \in \{0, 1\}^a$ (resp., $S_2(x, q) \in \{0, 1\}^a$).

On input $i \in [n]$, $r \in \{0, 1\}^L$, and answers $\alpha_1, \alpha_2 \in \{0, 1\}^a$, algorithm R outputs a bit $R(i, r, \alpha_1, \alpha_2)$, which is supposed to be a guess of the entry x_i .

The recovery condition: We denote by $\bar{A}(i, x)$ the random variable that represents the output of $R(i, r, S_1(x, q_1), S_2(x, q_2))$, where $(q_1, q_2) = Q(i, r)$ and the probability space is induced by the uniform distribution of $r \in \{0, 1\}^L$. Then, for every $i \in [n]$ and $x \in \{0, 1\}^n$, it must hold that $\Pr[\bar{A}(i, x) = x_i] \geq p$.

The secrecy condition: For $i \in [n]$, denote by $Q_1(i)$ (resp., $Q_2(i)$) the distribution induced on the first (resp., second) element of $Q(i, r)$ when r is uniformly distributed in $\{0, 1\}^L$. Then, for every $i, j \in [n]$, the distributions $Q_1(i)$ and $Q_1(j)$ (resp., $Q_2(i)$ and $Q_2(j)$) are δ -close (i.e., the statistical difference between them is at most δ).

Notice that we relax (and quantify) the security and recovery requirements; the traditional perfect requirements are obtained by setting $\delta = 0$ and $p = 1$. On the other hand, in the following, we restrict our attention to PIR schemes which have *linear answers*; that is, for every fixed query $q \in \{0, 1\}^t$, the servers' answers $S_1(x, q)$ and $S_2(x, q)$ are linear functions of x (each bit of $S_1(x, q)$ and each bit of $S_2(x, q)$ is a linear combination of the bits of x). The above-mentioned PIR scheme of Chor et. al. [4] satisfies this requirement.

Theorem 1.8 *Suppose there is a one-round, $(1 - \delta)$ -secure PIR scheme with 2 servers, linear answers, database size n , query size t , answer size a , and recovery probability $1/2 + \epsilon$. Suppose also that the user only uses k predetermined bits out of the a bits it receives as answer to each query. Then*

$$t > \frac{(\epsilon - \delta) \cdot n}{12 \cdot f(k, a)} - \log_2 f(k, a) - 3,$$

where $f(k, a) = \sum_{i=0}^k \binom{a}{i} \leq \min\{2^a, 2a^k\}$.

As immediate corollaries we conclude that

- Any (secure, one-round) 2-server PIR scheme with linear answers of constant length must have queries of linear (i.e., $\Omega(n)$) length. (This extends a simple lower bound (of $n - 1$ bits) on the length of queries in a 2-server PIR scheme with single-bit linear answers [4, Sec. 5.2].)
- Any (secure, one-round) 2-server PIR scheme with linear answers in which the user only uses one bit from each answer must have communication complexity $\Omega(\sqrt{n})$.
- Any (secure, one-round) 2-server PIR scheme with linear answers in which the user only uses k bits from each answer, k a constant, must have communication complexity $\Omega(n^{1/(k+1)})$.

In the abovementioned PIR scheme of Chor et. al. [4], both a and t are $O(n^{1/3})$, and $k = 4$. By a minor modification to that scheme, we can reduce k to 3. Thus the third lower bound asserts that for this case (i.e., $k = 3$), communication complexity of $\Omega(n^{1/4})$ is essential. We comment that the first two lower bounds are tight:

- There exists a (perfectly secure, one-round) 2-server PIR scheme that uses n -bit queries and linear answers that are single bits (cf., [4, Sec. 3.1]).
- There exists a (perfectly secure, one-round) 2-server, linear-answer PIR scheme in which the user uses only one bit from each \sqrt{n} bit-long answer, and the queries are also \sqrt{n} -bit long strings (e.g., by a minor modification of the scheme in [4, Sec. 3.2–3.3] as applied to $d = 2$).

Perspective: Computational security. We stress that the above results (as well as Section 5) refer to an information-theoretic notion of security. A relaxed notion of security, requiring only security with respect to polynomial-time servers, was put forward and first investigated by Chor and Gilboa [3]. Assuming the existence of one-way functions, for any $\epsilon > 0$, they presented 2-server computational-secure PIR schemes of communication complexity $O(n^\epsilon)$. Furthermore, their PIR schemes are one-round and use linear 1-bit answers. This stands in contrast to the lower bounds regarding the information-theoretic notion of security. Another PIR setting where computational security offers an advantage over information-theoretic security is the one of a single server (i.e., n bits is a lower bound in the case of information-theoretic security [4, Sec. 5.1], whereas communication complexity of $O(n^\epsilon)$ can be achieved for computationally-secure PIR's [7] under reasonable intractability assumptions).

1.3 Organization

Most of the paper is devoted to analysis of several types of locally decodable codes, and the application to private information retrieval is postponed to the last section (Section 5). Due to space considerations, two of the extensions mentioned above (i.e., to finite fields and block-block codes) are omitted. Full details can be found in our technical report [5].

2 Preliminaries

The notions and results in this section are mostly due to Katz and Trevisan [6]. In particular, their notion of smooth codes and its relation to locally decodable codes are central to our analysis. Here we generalize their definition to the case in which the message is over a non-Boolean alphabet.

2.1 Smooth Codes

Informally, a code is smooth if a corresponding local decoding algorithm “spreads its queries almost uniformly” (or, actually, does not query any code location too frequently).

Definition 2.1 For fixed c, ϵ , and integer q we say that $\mathbf{C} : \Sigma^n \rightarrow \Gamma^m$ is a (q, c, ϵ) -smooth code if there exists a probabilistic oracle machine A such that:

- In every invocation, A makes at most q queries non-adaptively.
- For every $x \in \{0, 1\}^n$ and for every $i \in [n]$, we have

$$\Pr[A^{\mathbf{C}(x)}(i) = x_i] \geq \frac{1}{2} + \epsilon.$$

- For every $i \in [n]$ and $j \in [m]$, the probability that on input i machine A queries index j is at most c/m .

(The probabilities are taken over the internal coin tosses of A .) An algorithm A satisfying the above requirements is called a (q, c, ϵ) -smooth decoding algorithm for \mathbf{C} .

We stress that the decoding condition in Definition 2.1 refers only to valid codewords, whereas the corresponding condition in Definition 1.1 refers to all oracles that are sufficiently close to valid codewords. To get a feeling for the smoothness condition note that if the decoding machine spreads its queries uniformly, then we would get $c = q$ (and this is the lowest possible value, assuming that the machine always makes q queries). It turns out that any locally decodable code is smooth, for suitable parameters and by possible modification of the decoding machine.

Theorem 2.2 (See Theorem 1 in [6]) Let $\mathbf{C} : \Sigma^n \rightarrow \Gamma^m$ be a (q, δ, ϵ) -locally decodable code. Then \mathbf{C} is also a $(q, q/\delta, \epsilon)$ -smooth code.

This is stated only for the case $\Sigma = \{0, 1\}$ in [6], but the proof applies to the general case as well.

2.2 The Recovery Graphs

Let $\mathbf{C} : \Sigma^n \rightarrow \Gamma^m$ be a $(2, c, \epsilon)$ -smooth code and let algorithm A be a (non-adaptive) $(2, c, \epsilon)$ -smooth decoding algorithm for \mathbf{C} . Let $\{q_1, q_2\}$ be a pair of elements of $[m]$. We say that a given invocation of A reads $\{q_1, q_2\}$ if the set of indices which A reads in that invocation is exactly $\{q_1, q_2\}$. We say that $\{q_1, q_2\}$ is good for i if:

$$\Pr[A^{\mathbf{C}(x)}(i) = x_i \mid A \text{ queries } \{q_1, q_2\}] > 1/2,$$

where the probability is taken over x uniformly chosen from $\{0, 1\}^n$, and over the internal coin tosses of A . For every $i \in [n]$, we consider the graph with edge set consisting of the set of good pairs.

Definition 2.3 Fixing a code $\mathbf{C} : \{0, 1\}^n \rightarrow \Gamma^m$ and a 2-query recovery algorithm A , the recovery graph for $i \in [n]$, denoted G_i , consists of the vertex set $[m]$ and the edge set E_i that equals the set of pairs $\{q_1, q_2\}$ that are good for i .

We have the following result about such graphs.

Lemma 2.4 ([6]) Let \mathbf{C} be a $(2, c, \epsilon)$ -smooth code and $\{G_i\}_{i=1}^n$ be the associated set of recovery graphs. Then, for every i , the graph $G_i = ([m], E_i)$ has a matching $M_i \subseteq E_i$ of size at least $\epsilon m/c$.

This is essentially Lemma 4 in [6], but, since we slightly changed the definition of the recovery graph (from [6]), and get slightly better bounds, we present a proof below.

Proof: We may assume without loss of generality that, for every $i \in [n]$ and $j_1, j_2 \in [m]$,

$$\Pr[A^{\mathbf{C}(x)}(i) = x_i \mid A \text{ queries } \{j_1, j_2\}] \geq \frac{1}{2} \quad (1)$$

where the probability is taken uniformly over $x \in \{0, 1\}^n$ and A 's internal coin tosses. (For example, we can modify A so that it outputs a random bit whenever $i \in [n]$ and $j_1, j_2 \in [m]$ do not satisfy Eq. (1).) Using a Markov argument, it follows that with probability at least 2ϵ , on input $i \in [n]$, algorithm A generates a pair that is good for i . In other words, with probability at least 2ϵ , the pair generated by $A(i)$ is an edge in G_i . Thus, if $C \subseteq [m]$ is a vertex cover of G_i , then the probability that $A(i)$ queries at least one element of C is at least 2ϵ . On the other hand, no element of $[m]$ is queried by A with probability greater than c/m , and so it follows that $|C| \geq (2\epsilon)/(c/m) = 2\epsilon m/c$. Since the size of the maximum matching in a graph is at least half the size of the minimum vertex cover, we conclude that G_i has a matching of size at least $\epsilon m/c$. ■

3 The Boolean Case – Proof of Theorem 1.3

3.1 Getting Rid Of Projected Bits

To simplify the rest of our analysis, we would like to get rid of bits in the range of the code that are identical to some input (data) bit. That is, we wish the code to be such that no single bit of the output is (always) equal to a particular bit of the input. We can accommodate this condition by fixing bits of the input that are identical to too many bits in the output. This gives the following lemma.

Lemma 3.1 For $n > 4c/\epsilon$, let $\mathbf{C} : \{0, 1\}^n \rightarrow \{0, 1\}^m$ be a (q, c, ϵ) -smooth code. Then there is another code $\mathbf{C}' : \{0, 1\}^{n'} \rightarrow \{0, 1\}^{m'}$ that has a $(q, c, \epsilon/2)$ -smooth reconstruction procedure A' , such that $n' \geq n/2$, $m' \leq m$, and for every i and j there exists an $x \in \{0, 1\}^{n'}$ such that the j th bit of $\mathbf{C}'(x)$ is different from x_i . Furthermore, if \mathbf{C} is a linear code, then so is \mathbf{C}' .

Thus lower bounds on the length of smooth codes satisfying the conclusion of the lemma yield lower bounds on general smooth codes.

Proof: Consider the set I of bits in the input that occur in more than a fraction $2/n$ of the bits of the output. Clearly, $|I| \leq n/2$. For each $i \in [n] \setminus I$, consider the behavior of the smooth reconstruction procedure $A^{C(x)}(i)$ for some x . Since $i \notin I$, at most a fraction $2/n$ of the bits of $C(x)$ contain copies of x_i . By the smoothness condition, such code bits are examined with probability at most $2c/n$, which is less than $\epsilon/2$ (provided that $n > 4c/\epsilon$). Thus, if we modify A such that it does not read such bits, we may decrease the probability that it recovers x_i by at most $\epsilon/2$, so the recovery condition is met.

We construct the code C' from C by omitting the output bits that are copies of any input bit $i \in [n]$, fixing arbitrary¹ values for the bits in I , “hardwiring” these values into C' , and modifying A so that it queries only bits in C' (rather than bits in C). Note that the fact that the length of C' may be shorter than the length of C only makes the smoothness condition easier to meet. ■

3.2 The Combinatorial Lemma

We will deal with the linear error-correcting code C' of Lemma 3.1. In the following we will use e_i to denote a vector in $\{0, 1\}^n$ that has 1 in the i -th coordinate and 0 elsewhere. We can identify our error-correcting code C' with a sequence of m' vectors $a_1, \dots, a_{m'} \in \{0, 1\}^{n'}$, such that the j th bit of $C(x)$ is $a_j \cdot x$. Recall that, by Lemma 3.1, none of these a_j 's equals any unit vector e_i . Let $\{G_i\}_{i=1}^{n'}$ be the sequence of recovery graphs associated with C' as in Lemma 2.4.

Lemma 3.2 *For every i , and for every $\{q_1, q_2\} \in E_i$, e_i is in the span of $\{a_{q_1}, a_{q_2}\}$.*

Proof: Suppose e_i is linearly independent of a_{q_1} and a_{q_2} . Then, for a random x , the value $x \cdot e_i$ is independent (in the statistical sense) of the values $x \cdot a_{q_1}$ and $x \cdot a_{q_2}$, and so it is not possible to gain any advantage in predicting x_i by looking at the q_1 -th and the q_2 -th bit of the encoding of x . ■

Since we are dealing with the field $\{0, 1\}$, when e_i is in the span of $\{a_{q_1}, a_{q_2}\}$ there are only three possibilities: either a_{q_1} or a_{q_2} equals e_i itself, or $e_i = a_{q_1} \oplus a_{q_2}$. But for C' (as in Lemma 3.1) the only possible case is that $e_i = a_{q_1} \oplus a_{q_2}$. Thus proving Theorem 1.3 reduces to proving the following result.

¹ Actually, in order to preserve linearity, these bits should all be set to zero. However, in fact, all our results apply also to affine codes.

Lemma 3.3 (Combinatorial Lemma) *Let a_1, \dots, a_m be elements of $\{0, 1\}^n$ such that for every $i \in [n]$ there is a set M_i of at least γm disjoint pairs of indices $\{j_1, j_2\}$ such that $e_i = a_{j_1} \oplus a_{j_2}$. Then $m \geq 2^{\gamma n}$. Furthermore, the conclusion holds even when the hypothesis only states that $\frac{1}{n} \sum_{i=1}^n |M_i| \geq \gamma m$.*

Below, we will present two alternative proofs of Lemma 3.3. Actually, the second proof yields a stronger lower-bound (of $m \geq 2^{2\gamma n}$, rather than $m \geq 2^{\gamma n}$). Combining all the above lemmas, we get:

Corollary 3.4 *Let $C : \{0, 1\}^n \rightarrow \{0, 1\}^m$ be a $(2, c, \epsilon)$ -smooth linear code. Then $m \geq 2^{\epsilon n/(4c)}$.*

Notice that Theorem 1.3 is an immediate consequence of Corollary 3.4 and Theorem 2.2.

Proof: We first apply Lemma 3.1 to obtain a $(2, c, \epsilon')$ -smooth linear code $C' : \{0, 1\}^{n'} \rightarrow \{0, 1\}^{m'}$, for $n' \geq n/2$, $m' \leq m$ and $\epsilon' = \epsilon/2$. Combining Lemmas 2.4 and 3.2, it follows that $\frac{1}{n'} \sum_{i=1}^{n'} |M_i| \geq \epsilon' m'/c$. Finally, applying Lemma 3.3, we get $m' \geq 2^{\epsilon' n'/c} \geq 2^{\epsilon n/(4c)}$, and using $m \geq m'$ the claim follows. ■

3.3 A Combinatorial Proof of Lemma 3.3

For starters, let us suppose that all the vectors a_1, \dots, a_m are different. In this special case, Lemma 3.3 is a consequence of the following known combinatorial result.²

Lemma 3.5 (See Appendix in [5]) *For any subset $S \subseteq \{0, 1\}^n$ of the hypercube, the number of edges of the hypercube having both endpoints in S is at most $\frac{1}{2}|S| \log_2 |S|$.*

Note that our (distinct) vectors a_1, \dots, a_m are all vertices of a hypercube, and we are assuming that, for every i , there are at least γm edges in the i th “direction” between such vertices. This gives a total of at least $\gamma m n$ edges, but this number has to be no more than $\frac{1}{2} m \log_2 m$, and so it follows that $m \geq 2^{2\gamma n}$.

To complete the proof of Lemma 3.3, we have to consider the case in which a_1, \dots, a_m are not all different. Note that an analogue of Lemma 3.5 does not hold in this case (e.g., if $a_1 = \dots = a_{m/2} = 0^n$ and $a_{(m/2)+1} = \dots = a_m = 10^{n-1}$ then we get $(m/2)^2$ edges).³

² The proof of Theorem 2 in [2, Sec. 16] implies that the subset $S \subseteq \{0, 1\}^n$ of given size m for which the number of internal edges is maximum is the set of the first $m = |S|$ strings in lexicographic order (of $\{0, 1\}^n$). Since each such vertex has at most $\lceil \log_2 m \rceil$ internal edges, we get an upper-bound of $\frac{1}{2}|S| \lceil \log_2 |S| \rceil$ on the number of internal edges. Indeed, the difference is of little significance in the context of our work.

³ Note that this example does not violate Lemma 3.3: for every sequence of M_i 's as in Lemma 3.3, it holds that $\sum_{i=1}^n |M_i| \leq 1$ (since $|M_1| \leq 1$ and all the other M_i 's must be empty). Thus, the “furthermore hypothesis” only holds with $\gamma \leq 1/(nm)$, implying a lower bound of $m \geq 2^{\gamma n} \geq 2$ (which indeed holds).

For every $a \in \{0, 1\}^n$, let us denote by ν_a the number of indices j such that $a_j = a$ (so that $\sum_{a \in \{0, 1\}^n} \nu_a = m$). That is, ν_a is the multiplicity of the vector a in the sequence a_1, \dots, a_m . For every k , let us denote by S_k the set of vectors a such that $\nu_a \geq k$, and let $s_k = |S_k|$; observe that

$$\sum_k s_k = m, \quad (2)$$

because each vector a that occurs in the sequence a_1, \dots, a_m is counted exactly ν_a times. Finally, define $\chi(a, j)$ to be 1 if $\nu_a \geq j$ and to be 0 otherwise. With this new piece of notation we can write

$$\sum_{a \in \{0, 1\}^n} \sum_{k \geq 1} \chi(a, k) = m, \quad (3)$$

and we also note that for any two vectors $a, b \in \{0, 1\}^n$, we have

$$\min\{\nu_a, \nu_b\} = \sum_{k \geq 1} \chi(\nu_a, k) \chi(\nu_b, k). \quad (4)$$

Now we would like to argue that for every i , the following upper bound holds on the size of the matching M_i :

$$|M_i| \leq \sum_{a, b: a \oplus b = e_i} \min\{\nu_a, \nu_b\}. \quad (5)$$

Indeed, for starters we have by definition that M_i is the set of all pairs $\{j_1, j_2\}$ such that $a_{j_1} \oplus a_{j_2} = e_i$, and that all such pairs are disjoint. Let us fix two vectors a and b such that $a \oplus b = e_i$, and consider how many possible pairs $\{j_1, j_2\}$ can belong to M_i subject to $a_{j_1} = a$ and $a_{j_2} = b$; since the pairs have to be disjoint, both ν_a and ν_b are upper bounds on the number of such possible pairs. Summing over all choices of a and b gives the bound of (5).

Combining the lemma's hypothesis with Equations (5) and (4), we get

$$\begin{aligned} \gamma mn &\leq \sum_{i=1}^n |M_i| \\ &\leq \sum_{i=1}^n \sum_{a \in \{0, 1\}^n} \min\{\nu_a, \nu_{a \oplus e_i}\} \\ &= \sum_{i=1}^n \sum_{a \in \{0, 1\}^n} \sum_{k \geq 1} \chi(\nu_a, k) \chi(\nu_{a \oplus e_i}, k) \end{aligned}$$

and so

$$\gamma mn \leq \sum_{k \geq 1} \sum_{i=1}^n \sum_{a \in \{0, 1\}^n} \chi(\nu_a, k) \chi(\nu_{a \oplus e_i}, k). \quad (6)$$

Note that $\sum_{i=1}^n \sum_{a \in \{0, 1\}^n} \chi(\nu_a, k) \chi(\nu_{a \oplus e_i}, k)$ counts (twice) the number of hypercube edges with both endpoints in S_k . Thus, by Lemma 3.5, we have, for every k , that

$$\begin{aligned} \sum_{i=1}^n \sum_{a \in \{0, 1\}^n} \chi(\nu_a, k) \chi(\nu_{a \oplus e_i}, k) &\leq 2 \cdot \frac{1}{2} |S_k| \log_2 |S_k| \\ &= s_k \log_2 s_k \\ &\leq s_k \cdot \log_2 m. \end{aligned}$$

Combining this inequality with (6), and recalling (2), we have

$$\gamma mn \leq \sum_k s_k \cdot \log_2 m = m \cdot \log_2 m,$$

from which it follows that $m \geq 2^{\gamma n}$.

3.4 An Alternative Proof of Lemma 3.3

The ‘‘information-theoretic’’ proof in this section is due to Alex Samorodnitsky, and was suggested to us after we found the combinatorial proof presented in the previous subsection.

Let X be a random variable uniformly distributed in the multiset $\{a_1, \dots, a_m\}$. We will write $X = X_1 X_2 \dots X_n$, where X_i denotes the i th bit of X , and $X_{i,j}$ denotes $X_i \dots X_j$. We consider the entropy of X , denoted $H(X)$. On one hand, $H(X) \leq \log_2 m$. On the other hand, we will prove that $H(X) \geq 2\gamma n$, and Lemma 3.3 will follow immediately.

We can express the entropy of X as

$$\begin{aligned} H(X) &= H(X_1) + H(X_2|X_1) \\ &\quad + \dots + H(X_n|X_1 \dots X_{n-1}). \end{aligned}$$

The value of the i th term, $H(X_i|X_1 \dots X_{i-1}) = H(X_i|X_{1,i-1})$, is given by the following formula:

$$\begin{aligned} H(X_i|X_{1,i-1}) &= \sum_{b \in \{0, 1\}^{i-1}} \Pr[X_{1,i-1} = b] \cdot H(X_i|X_{1,i-1} = b). \quad (7) \end{aligned}$$

Observe that for any 0-1 random variable Y (in our case $Y = (X_i|X_{1,i-1} = b)$), with $p \stackrel{\text{def}}{=} \Pr(Y = 1)$, we have $H(Y) = H_2(p)$, where $H_2(x) = x \log_2(1/x) + (1-x) \log_2(1/(1-x)) \geq 2 \cdot \min(x, 1-x)$ is the binary entropy function.⁴ So Eq. (7) is at least

$$\begin{aligned} &\sum_{b \in \{0, 1\}^{i-1}} \Pr[X_{1,i-1} = b] \\ &\quad \cdot 2 \cdot \min_{\sigma \in \{0, 1\}} \{\Pr[X_i = \sigma | X_{1,i-1} = b]\}. \quad (8) \end{aligned}$$

⁴We claim that, for $x \in [0, 0.5]$, it holds that $H_2(x) \geq 2x$ (whereas a bound of $H_2(x) \geq x$ is obvious). The claim can be verified by noting that $f(x) \stackrel{\text{def}}{=} H_2(x) - 2x$ is convex in that interval, and that $f(0) = 0 = f(1/2)$.

Now, under any conditioning, the probability that X is an endpoint of an edge in M_i equals the sum over $\sigma \in \{0, 1\}$ of the probabilities that $X_i = \sigma$ and X is an endpoint of an edge in the matching M_i (which matches events of the type $X_i = 0$ with events of the type $X_i = 1$). Thus, each of the two probabilities in the sum is bounded above by $\min(\Pr[X_i = 0|\text{cond}], \Pr[X_i = 1|\text{cond}])$. Thus, $\Pr[X \text{ is an endpoint of } e \in M_i | X_{1,i-1} = b]$ is bounded above by $2 \cdot \min\{\Pr[X_i = 0 | X_{1,i-1} = b], \Pr[X_i = 1 | X_{1,i-1} = b]\}$, and Eq. (8) is bounded below by

$$\begin{aligned} & \sum_{b \in \{0,1\}^{i-1}} \Pr[X_{1,i-1} = b] \\ & \cdot \Pr[X \text{ is an endpoint in an edge of } M_i | X_{1,i-1} = b] \\ & = \Pr[X \text{ is an endpoint in an edge of } M_i] \\ & = \frac{2|M_i|}{m} \geq 2\gamma. \end{aligned}$$

Then $H(X) \geq 2\gamma n$ and so $m \geq 2^{2\gamma n}$.

Comment: Note that the lower bound established here (i.e., $m \geq 2^{2\gamma n}$) is a square of the lower-bound claimed in Lemma 3.3. Furthermore, this stronger lower-bound is tight, and implies Lemma 3.5 as a special case.⁵

4 Extension To Binary Linear Block Codes – The Proof of Theorem 1.5

In this section we deal with linear codes mapping $\{0, 1\}^n$ to $(\{0, 1\}^\ell)^m$, where the case $\ell = 1$ corresponds to the main result (presented in Section 3). Thus each output symbol is an ℓ -bit long string, where each of these bits is a linear combination of the n input bits. We show that providing lower bounds for the general case reduces to providing lower bounds for the special case of $\ell = 1$.

4.1 Reduction to the Boolean case

Lemma 4.1 *Let $\mathbf{C} : \{0, 1\}^n \rightarrow (\{0, 1\}^\ell)^m$ be a (q, c, ϵ) -smooth linear error-correcting code. Then there is a code $\mathbf{C}' : \{0, 1\}^n \rightarrow \{0, 1\}^{2^{\ell \cdot m}}$ that is $(q, c \cdot 2^\ell, \epsilon)$ -smooth. Furthermore, suppose that \mathbf{C} has a decoding algorithm that uses only k predetermined bits out of the ℓ bits that it receives as answer to each query. Then there is a code $\mathbf{C}'' : \{0, 1\}^n \rightarrow \{0, 1\}^{t \cdot m}$ that is $(q, c \cdot t, \epsilon)$ -smooth, where $t = \sum_{i=0}^k \binom{\ell}{i}$.*

⁵ Specifically, the set of edges $E(S, S)$ with both endpoints in S can be partitioned into matchings M_i 's as in Lemma 3.3. Letting $\gamma = (\sum_i |M_i|)/(n|S|)$, and applying the stronger bound (for Lemma 3.3), we get $|S| \geq 2^{2\gamma n} = 2^{2 \sum_i |M_i|/|S|}$. Thus, $\log_2 |S| \geq 2|E(S, S)|/|S|$, which implies $|E(S, S)| \leq (1/2)|S| \log_2 |S|$.

Proof: Let $x \in \{0, 1\}^n$. We define $\mathbf{C}'(x)$ as follows: for every $j \in [m]$ and for every $a \in \{0, 1\}^\ell$, the entry of $\mathbf{C}'(x)$ indexed by (j, a) contains the inner product between the j th (ℓ -bit long) block of $\mathbf{C}(x)$ and the (ℓ -bit long) string a . This encoding has length $m' \stackrel{\text{def}}{=} 2^\ell m$. We now describe a smooth decoding procedure for \mathbf{C}' .

Let A be the $(2, c, \epsilon)$ -smooth decoding procedure for \mathbf{C} . The smooth decoding procedure A' for \mathbf{C}' will first simulate A , and get two queries (j_1, j_2) . If x_i is in the span of $\mathbf{C}(x)_{j_1}$ and $\mathbf{C}(x)_{j_2}$, then A' will reconstruct x_i as a linear combination of $\mathbf{C}(x)_{j_1}$ and $\mathbf{C}(x)_{j_2}$, a computation that can be done by looking at two entries of $\mathbf{C}'(x)$ (i.e., specifically the entries (j_1, a_1) and (j_2, a_2) , where $x_i = \langle a_1, \mathbf{C}(x)_{j_1} \rangle + \langle a_2, \mathbf{C}(x)_{j_2} \rangle$). If x_i is not in the span of $\mathbf{C}(x)_{j_1}$ and $\mathbf{C}(x)_{j_2}$, then A' will output a random guess. As argued in the proof of Lemma 2.4, with probability at least 2ϵ , algorithm A (on input i) samples a pair (j_1, j_2) that is good for i (i.e., allows reconstruction with average success probability above $1/2$, when averaging over all possible x 's). However, whenever (j_1, j_2) is good for i , it must be the case that x_i is in the span of $\mathbf{C}(x)_{j_1}$ and $\mathbf{C}(x)_{j_2}$, and A' correctly reconstructs x_i . Combining these two observations, we bound the reconstruction probability of A' below by $2\epsilon \cdot 1 + (1 - 2\epsilon) \cdot (1/2) = 1/2 + \epsilon$ (as required). Turning to the smoothness condition, observe that each entry in $\mathbf{C}'(x)$ is queried with probability at most c/m , which equals $(2^\ell \cdot c)/m'$ as required.

In order to prove the “furthermore” part, we do a similar construction, except that the entries of $\mathbf{C}''(x)$ correspond to pairs (j, a) where $j \in [m]$ and $a \in \{0, 1\}^n$ is a vector of weight at most k . When introducing the decoding procedure A'' (for \mathbf{C}''), we refer not only to the queries made by A but also the the predetermined bit locations in the answer that are inspected by A . Specifically, A'' first simulates A , and gets two queries (j_1, j_2) as well as two corresponding sets of bit locations $S_1, S_2 \subseteq [\ell]$. If x_i is in the span of the bit positions S_1 in $\mathbf{C}(x)_{j_1}$ and the bit positions S_2 in $\mathbf{C}(x)_{j_2}$, then A'' will reconstruct x_i as a linear combination of these bit positions, a computation that can be done by looking at two entries of $\mathbf{C}''(x)$, since $|S_1|, |S_2| \leq k$. In the analysis we note that whenever a pair of queries (made by A) is good for i , it must be the case that x_i is in the span of the bits of $\mathbf{C}(x)_{j_1}$ and $\mathbf{C}(x)_{j_2}$ that are inspected by A , and A'' correctly reconstructs x_i . ■

4.2 Consequences

Combining Lemma 4.1 and Corollary 3.4, we obtain the following result.

Corollary 4.2 *Let $\mathbf{C} : \{0, 1\}^n \rightarrow (\{0, 1\}^\ell)^m$ be a (q, c, ϵ) -smooth linear error-correcting code. Then $m \geq (1/2^\ell) \cdot 2^{\epsilon n/4 \cdot 2^{\ell \cdot c}}$. Furthermore, if \mathbf{C} has a decoding algorithm that*

uses only k of the ℓ bits that it receives as answer to each query, then $m \geq (1/t) \cdot 2^{\epsilon n/4 \cdot t \cdot c}$, where $t = \sum_{i=0}^k \binom{\ell}{i}$.

Theorem 1.5 follows by combining Corollary 4.2 and Theorem 2.2.

5 Lower Bounds For Private Information Retrieval – Proof of Theorem 1.8

The main result of this section is a reduction showing that a one-round PIR system can be converted into a smooth error-correcting code. This transformation preserves linearity, and hence, combined with the lower bound for smooth linear codes, yields a lower bound for linear one-round PIR systems.

5.1 Constructing Smooth Codes Based on PIR Schemes

Actually, we consider a relaxed notion of a PIR. First, recovery is not required to always be correct but rather only to be correct with probability at least $1/2 + \epsilon$, where the probability is taken over the PIR's randomization for any fixed input (i.e., a database and a desired bit). Second, we do not require perfect secrecy (i.e., $\delta = 0$), but rather that the distributions of each query for each desired bit are at pairwise statistical distance at most δ .

Lemma 5.1 *Suppose there is a one-round, $(1 - \delta)$ -secure PIR scheme with two servers, database size n , query size t , answer size a , and recovery probability at least $1/2 + \epsilon$. Then there is a $(2, 3, \epsilon - \delta)$ -smooth error-correcting code $\mathbf{C}: \{0, 1\}^n \rightarrow (\{0, 1\}^a)^m$, where $m \leq 6 \cdot 2^t$. Furthermore:*

1. *If in the PIR scheme the answer bits are a linear combination of the data, then \mathbf{C} is linear.*
2. *If, in the PIR scheme, the user only uses k predetermined bits out of the a bits it receives as an answer to each question, then the same property is true for the decoding algorithm of \mathbf{C} .*

Proof: Let us first develop some intuition about the proof. By enumerating all possible answers from either server, we can view the PIR system as encoding the database $x \in \{0, 1\}^n$ as a string $PIR(x) \in (\{0, 1\}^a)^l$, where $l = 2 \cdot 2^t$. The user can reconstruct one bit x_i of the database with advantage ϵ by looking at two entries of the encoded string $PIR(x)$. For any i and j , the distribution of the first entry read into $PIR(x)$ when reconstructing x_i is δ -close to the distribution of the first entry read into $PIR(x)$ when reconstructing x_j (and similarly for the second entry). Instead of this closeness property, we would like to have a smoothness property, that is, we would like each entry to be read with low probability. We are willing to make the

encoding be slightly longer in order to achieve this goal. We will achieve this goal by duplicating entries that have a high probability of being read.

Suppose, to start, that $\delta = 0$. Then, for every j , the probability that entry j is queried by the reconstruction algorithm (as a first query or as a second query) is a fixed value p_j (independent of which bit of the database the user wants to reconstruct); note that $\sum_j p_j = 2$. We will replicate entry j of the encoding $n_j = \lceil p_j \cdot l \rceil$ times, denoting by $\mathbf{C}(x)$ this new encoding (with repetitions) of x . Recall that $PIR(x) \in (\{0, 1\}^a)^l$ (and we will show that $\mathbf{C}(x) \in (\{0, 1\}^a)^{O(l)}$).

A reconstruction algorithm for x_i from $\mathbf{C}(x)$ will generate queries j_1, j_2 as in the reconstruction algorithm that accesses $PIR(x)$. The algorithm then picks at random one of the n_{j_1} copies of the j_1 th entry and one of the n_{j_2} copies of the j_2 th entry, and then accesses these selected two entries in $\mathbf{C}(x)$. Clearly, the advantage in decoding x_i remains the same. Regarding smoothness, let us consider an entry j in $PIR(x)$. If $p_j \leq 1/l$, then the corresponding (unique) bit in $C(x)$ is accessed with probability $p_j \leq 1/l$. Otherwise (i.e., $p_j > 1/l$), the j th entry is replicated $n_j = \lceil p_j l \rceil > 1$ times, and each copy is accessed with probability p_j/n_j , which is

$$\frac{p_j}{\lceil p_j l \rceil} \leq \frac{p_j}{p_j l} = \frac{1}{l}.$$

The length of the new encoding is $m = \sum_{j=1}^l n_j$, and we have

$$\begin{aligned} m &= \sum_{j:p_j \leq 1/l} \lceil p_j l \rceil + \sum_{j:p_j > 1/l} \lceil p_j l \rceil \\ &\leq \sum_{j:p_j \leq 1/l} 1 + \sum_{j:p_j > 1/l} (1 + p_j l) \\ &\leq l + \sum_j p_j l \\ &= 3l = 6 \cdot 2^t. \end{aligned}$$

Recall that no entry is queried with probability higher than $1/l$, which (using $m \leq 3l$) is bounded above by $3/m$.

Consider now the general case in which the query distributions for x_{i_1} and x_{i_2} are only guaranteed to be δ -close. We apply the previously described construction using the distribution of queries for x_1 . When we want to reconstruct x_i we proceed as follows. For every j , let p_j be the probability that j is queried when reconstructing x_1 and let q_j be the probability that j is queried when reconstructing x_i . Note that $\sum_j p_j = \sum_j q_j = 2$ and that $\sum_j |p_j - q_j| \leq 4\delta$, and so $\sum_{j:q_j > p_j} (q_j - p_j) \leq 2\delta$. We sample queries j_1, j_2 as in the original algorithm for x_i (modified so as to choose a random copy, if the required entry has multiple copies), and then if $q_{j_1} \leq p_{j_1}$, we proceed to make query j_1 . If $q_{j_1} > p_{j_1}$, then we read query j_1 with probability p_{j_1}/q_{j_1} and we enter a

“failure mode” with the remaining probability. In failure mode, bit x_i is just guessed randomly. Query j_2 is handled similarly.

Observe that the smoothness requirement is satisfied as before (since each bit corresponding to the original query j is accessed with probability $\min\{q_j, p_j\}/n_j \leq p_j/n_j \leq 1/l$). The probability of entering the failure mode is $\sum_{j:q_j>p_j} (q_j - p_j) \leq 2\delta$, and when the failure mode is entered, the probability of guessing x_i correctly is exactly one half. Thus, in the worst case, failures subtract δ of the probability of guessing x_i correctly, and so the overall probability of guessing x_i right is at least $1/2 + \epsilon - \delta$. ■

5.2 Consequences

Theorem 1.8 follows by combining Lemma 5.1 and Corollary 4.2. Specifically, using $m \leq 6 \cdot 2^t$, a smoothness bound of $c = 3$ and recovery advantage $\epsilon - \delta$, we have $6 \cdot 2^t \geq \frac{1}{f(k,a)} \cdot 2^{\frac{(\epsilon-\delta) \cdot n}{4^{3-f(k,a)}}}$, and Theorem 1.8 follows.

Acknowledgments

We are grateful to Alex Samorodnitsky for suggesting to us the information-theoretic proof of Lemma 3.3 and allowing us to present it in Section 3.4. Thanks also to Noga Alon for providing us with a proof of Lemma 3.5 and allowing us to reproduce it in our technical report [5].

References

- [1] A. Ambainis. An Upper Bound On The Communication Complexity of Private Information Retrieval. In *24th ICALP*, Springer, Lecture Notes in Computer Science, Vol. 1256, pages 401–407, 1997.
- [2] B. Bollobás. *Combinatorics*. Cambridge University Press, 1986.
- [3] B. Chor and N. Gilboa. Computationally-Private Information Retrieval. In *29th STOC*, pp. 304–313, 1997.
- [4] B. Chor, O. Goldreich, E. Kushilevitz and M. Sudan, Private Information Retrieval. *Journal of the ACM*, Vol. 45, No. 6, pages 965–982, November 1998.
- [5] O. Goldreich, H. Karloff, L.J. Schulman and L. Trevisan, Lower Bounds for Linear Locally Decodable Codes and Private Information Retrieval. *ECCC*, TR01-080, 2001.
- [6] J. Katz and L. Trevisan. On The Efficiency Of Local Decoding Procedures For Error-Correcting Codes. In *32nd STOC*, 2000.
- [7] E. Kushilevitz and R. Ostrovsky. Replication is Not Needed: Single Database, Computationally-Private Information Retrieval. In *38th FOCS*, pages 364–373, 1997.