

Lecture 13

*Lecture date: February 26, 2024**Scribes: Benjamin Gillmore*

1 Private information retrieval (PIR)

1.1 Overview

Consider a game with user U and a database DB . The database is, from the theoretical perspective, just an array of bits from 1 to $\{1..m\}$. What user wants to do is retrieve location i of the database, the element x_i . Furthermore, U doesn't want the DB to know location i . This is comparable in concept to oblivious transfer where U wishes to hide which bit he aims to retrieve.

Definition 1 (Communication Complexity (CC)) *Given a protocol between two players, communication complexity is the worst case number of bits that have to be transferred to achieve their goal.*

What makes PIR non-trivial is that we want the CC of PIR, both from U and DB , to be strictly less than n . Should this not be the case, consider the trivial case where U asks for the entire database and does not disclose which bit they are interested in.

It is also a bad solution to say for in addition to i , U will ask for 300 other locations. While DB doesn't know which out of the 301 locations is the location U is interested in, the DB can eliminate some locations that U is not interested in. That is, this approach still reveals information about i . Given the standard notion of indistinguishably, it should be that for any location accessed, the DB cannot rule it out. Furthermore, after the communication is over, the DB 's ability to decide if U is interested in location i or location j is the same.

1.2 Protocol 1

Private information retrieval was introduced in a paper by Chor, Goldreich, Kushilevitz, Sudan (96). They basically said two things. First they proved that this game is theoretically impossible with one database. However, if you have two identical databases DB_1 and DB_2 that cannot communicate, that is that DB_1 cannot communicate U 's interactions with DB_2 and likewise DB_2 to DB_1 , then this notion of PIR is possible.

Consider a database that is n bits, form DB_1 and DB_2 such that they each form identical matrices of size $\sqrt{n} * \sqrt{n}$ and cannot collude against U . U wants to know some position i in the database but since the database has been rearranged, the location U is interested in is at some new position, say (i, j) . The user can pick some random bit vector r of size \sqrt{n} and sends r to DB_1 . DB_1 returns a bit wise exclusive or of all columns where $r = 1$. Next U sends DB_2 a bit vector r_j which is exactly the vector r with the j 'th bit flipped. DB_2 does the exact same thing as DB_1 and returns to U the exclusive or of the columns referenced by r_j . It is trivial to see that if U takes these two vectors and exclusive or them together then they receive the entirety of j^* , or the j 'th column of the database. U can simply look at the i 'th bit of j^* to find the element at (i, j) . Consider the CC of this protocol, U sends a vector of size \sqrt{n} to both of the databases and the databases each send a vector of size \sqrt{n} back, which results in a total CC of $4\sqrt{n}$.

1.3 Matching Vector Family

Note: Please see comment in tex.

A matching vector family are n vectors $\{v_1, v_2, \dots, v_n\}$ and $\{w_1, w_2, \dots, w_n\}$ that have the following property $v_i \bullet w_i = 0 : \forall i$ and $v_i \bullet w_j \neq 0 : \forall i, j (i \neq j)$ A paper by Bhowmick, Dvir, and Lovett showed how to construct a matching vector family where the length of this vector is of n^ϵ . They also proved that if can reduce the size of this vector to poly-logarithmic then you get a big result. However, this remains a big open question in combinatorics.

1.4 Protocol 2

If you have multiple non-colluding databases can you do better? The answer is yes, sort of. Suppose you have four identical non-communicating databases DB_1, DB_2, DB_3, DB_4 organized as a matrix of size $\sqrt{n} * \sqrt{n}$ just as before. Again, U is interested in some location (i, j) . U sends a random vectors x and y of size \sqrt{n} bits to DB_1 . DB_1 searches the database for locations where $DB[i^*, j^*] = 1 : \forall i^*, j^* (x_{i^*} = x_{j^*} = 1)$, exclusive or's the results together and returns this single bit sum to U . Next, U repeats this process with DB_2 but uses x_i and y where x_i is x with the i 'th bit flipped. Again, U repeats this process with DB_3 but uses x and y_j where y_j is y with the j 'th bit flipped. Finally, U repeats this process with DB_4 but uses x_i and y_j . Let the results from each database be s_1, s_2, s_3, s_4 and U can now calculate $DB[i, j] = s_1 \otimes s_2 \otimes s_3 \otimes s_4$.

To illustrate this further consider an index where $x_{i^*} = 1, y_{j^*} = 1$ and $i^* \neq i, j^* \neq j$. We know every database will return a sum with this bit included, so U will effectively be calculating $DB_1[i^*, j^*] \otimes DB_2[i^*, j^*] \otimes DB_3[i^*, j^*] \otimes DB_4[i^*, j^*]$. Since all databases are identical, we can infer all the results to be identical, and also know a sum of any even number of identical elements mod 2 to be 0. So for all rows and columns except for i and

j , everything appears four times and everything gets cancelled.

In the case of DB_2 and DB_3 , where exclusively either the i 'th bit of x or the j 'th bit of y is referenced by U , we can infer that the databases will sum either the inner product of i 'th row, and the j 'th column or the inner product of i 'th row and the j 'th column in their respective s , but not the bit at $DB[x_i, y_j]$. Thus, the elements referenced in these cases are factored into the final sum once here. Finally, the case of DB_4 where both row i and column j are referenced by U , it is trivial to see that the element at $DB[x_i, y_j]$ will be included in the sum s_4 . Also notice that the inner product of all other bits in the i 'th row, and the j 'th column will be summed in s_4 , these will cancel out the the results from s_2 and s_3 and leave U with only the element at index (i, j) . As a side note, the i 'th bit of x and the j 'th bit of y need not be originally 0, in fact they must be random to hide U 's intentions to the database.

Notice, the total CC here is of $8\sqrt{n} + 4$ bits with $8\sqrt{n}$ bits U must transmit and only 4 bits that U must receive.

1.5 More Dimensions

While this looks like fun in two dimensions, consider a 3 dimensional case. Suppose you have 8 identical non-colluding databases $\{DB_1, DB_1, \dots, DB_8\}$ and random vectors x, y , and z of length $n^{1/3}$. Similarly to before, the databases can pack their contents into a 3-dimensional cube which the the user U can reference with the indices i, j , and k . Again, U creates vectors x^*, y^* , and z^* with the i 'th, j 'th, and k 'th bit flipped respectively. U now sends $DB_1(x_i, x_j, x_k)$ and like before returns the sum mod 2 of the referenced bits. U repeats the process with DB_2 with (x_{i^*}, y_j, z_k) , DB_3 with (x_i, y_{j^*}, z_k) , DB_4 with (x_{i^*}, y_{j^*}, z_k) , DB_5 with (x_i, y_j, z_{k^*}) , DB_6 with (x_{i^*}, y_j, z_{k^*}) , DB_7 with (x_i, y_{j^*}, z_{k^*}) , and DB_8 with $(x_{i^*}, y_{j^*}, z_{k^*})$. Let the results from each database be the respective sums $\{s_1, s_2, \dots, s_8\}$. The claim is that U can now calculate $DB[i, j, k] = s_1 \otimes s_2 \otimes \dots \otimes s_8$. This can be proven similarly to before by counting which bits appear an odd number of times and which appear an even number of times. Notice that the total CC here is $16n^{1/3} + 8$ bits, with U transmitting $16n^{1/3}$ bits and receiving 8 bits.

This seems to be utterly ridiculous because now we need to have to have 8 non-colluding databases. Chor, Goldreich, Sudan observed that we can actually simulate all 8 databases with only 2 databases. The idea is as follows U only gives x, y , and z to DB_1 . Notice that from before the CC from U to DB_1 is unbalanced in that U sends $DB_1 \sqrt{n}$ bits and DB_1 sends U a single bit. So what DB_1 does instead is iterate all possible i bit flips (x^*) of x and return to U the $n^{1/3}$ inner products of all possible x^*, y , and z . U and DB_1 now repeat this process for all possible j bit flips (y^*) of y and all possible k bit flips (z^*) of z . U now repeats this process with DB_2 with the i 'th, j 'th, and k 'th bit flipped of x, y , and z respectively.

Note: Please see comment in tex.

Goldreich, Karloff, Schulman, and Trevisan showed that if you do exclusive or summing, $n^{1/3}$ is a lower bound for protocols such as the above, but it turns out if you use more sophisticated math you can actually go all the way to n^ϵ and it is expected that there exists a poly-logarithmic CC solution for cases with only 2 or 3 databases.

2 Locally Decodable Codes (LDC)

The notion of a locally decodable code is that you don't have to read the actual code word of interest, rather one can read just a few locations that are correlated with the code word to recover its value. Katz and Trevisan observed the following connection between PIR and LDC. Consider any database PIR uniform random solution, say with two databases, U asks some question x_i to DB_1 and some question y_j to DB_2 , as long as the databases do not compare notes, these questions look completely random, and do not reveal the location U is interested in. Let's imagine the following absolutely preposterous error correcting code, lets enumerate all possible 2^n questions for x_i from 000...0 to 111...1. Similarly lets do the same for DB_2 , enumerate all possible questions for y_j . Now we can say that our encoding for x is this vector of all possible results from the two databases. Consider an adversary that has the ability to corrupt $1/8$ 'th of the bits in each of the databases. While, the question x that the U asks to DB_1 is uniform random so the probability U is receiving non-corrupted data is $7/8$ from each database individually. Together the probability that U receives corrupted data is $1/4$ so therefore with probability $3/4$ U will be getting a correct answer. However, U need not read the correct code word, they may choose to any other two instances and recover the location of interest. Thus any PIR scheme where the questions are uniform random implies LDC. As it turns out LDC implies PIR as well.

3 PIR with Homomorphic Encryption

Imagine that we have probabilistic homomorphic encryption as such $E(x)*E(y) = E(x+y)$. There is a single database DB and there is a user U what wants to know location i of the database. DB will partition the database into blocks of size \sqrt{n} and U will send to the database an encryption of $E(0)$ and $E(1)$ for each block. DB will replace every entry of 0 with $E(0)$ and every entry of 1 with $E(1)$. However, in actuality U lies to DB . U sends the DB $E(0)$ and $E(1)$ for blocks they are interested in and two $E(0)$'s for blocks that they are not interested in. Next, the database homomorphically multiplies every i 'th entry of every block B_i together and returns the result to U . Because U did not tell DB which B_i they lied about and only told the truth for the B_i they are interested in, DB has no notion of which B_i is of interest to U . Yet still DB is returning to U an encryption of the

B_i of interest because DB is unknowingly performing either $E(0) * E(0) = E(0 + 0)$ or $E(0) * E(B_i) = E(0 + B_i)$.

4 PIR and 1-out-of-N OT

Note: Again, please see latex comments for attribution.

Suppose we have a database DB and a user U . U wants to retrieve an element i without showing DB which i is being retrieved and DB wants to make sure U gets i and nothing else. Moni Naor's idea was to have DB create 2 groupings of $\log(n)$ keys of $\log(n)$ bits in length, k_{i0} and k_{i1} . For each entry, the database will encrypt j 'th bit of each entry i with k_{i0} if $k_{ij} = 0$ and k_{i1} if $k_{ij} = 1$. For each of the $\log(n)$ key pairs, U now uses 1-out-of-2-OT to get one key out and selects the address they are interested in. Now the user can play the standard PIR where all the bits are encrypted using these different combination of keys, receives the encrypted results, and decrypts.

5 Doubly Efficient PIR

Boneh and Corrigan-Gibbs came up with the notion of Doubly Efficient PIR. Suppose we have 2 databases DB_1 , DB_2 and a user U in which there is an online phase and an offline phase. In the online phase, U chooses a random subset of \sqrt{n} bits and the DB_1 exclusive or's the bits, returns this result to the user, and they repeat this process \sqrt{n} times. This can be done efficiently by specifying a seed of a PRF. So U knows \sqrt{n} sums of random subsets of bits. Now in the offline phase, U wants to know location i , he also knows that location i was part of the sum in the \sqrt{n} subsets they asked DB_1 for. What U can do is tell DB_2 some of the same random subsets he had told to DB_1 . DB_2 now repeats the process DB_1 did before and returns their sums to U . U and DB_2 repeat this process a few times until U can construct a sum of all subsets he had DB_1 sum without including the results that contained location i . All U must do now is subtract the sum they constructed from DB_2 's result from DB_1 's result, or exclusive or them together, to retrieve location i .