

Efficient Private Matching and Set Intersection

Mike Freedman, NYU
Kobbi Nissim, MSR
Benny Pinkas, HP Labs

(To appear in EUROCRYPT 2004)

A Story...



We think patients are misusing prescriptions to obtain drugs...

Here too..

We could share our lists of patients?

But, what about HIPAA? And we're competitors!

Have you heard of "secure function evaluation" ?

This is all "theory". It can't be efficient.

A Story...



1.Improvements to generic primitives (SFE, OT)

2.Improvements in specific protocol examples

We could share our lists of patients?

Have you heard of "secure function evaluation" ?

This is all "theory". It can't be efficient.

The Scenario



Client



Server

Input: $X = x_1 \dots x_k$

$Y = y_1 \dots y_k$

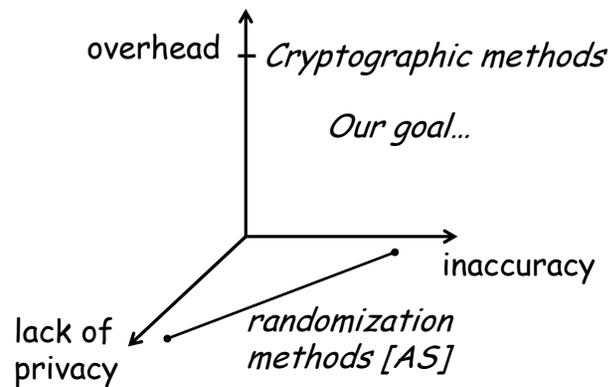
Output: $X \cap Y$ only

nothing

- Enterprises and government holding sensitive databases
- Peer-to-Peer networks
- Mobile wireless crowds (PDAs, cell phones)

Credit rating, CAPS II, shared interests (research, music), genetic compatibility, etc

Crypto vs. randomization methods



Related work

- Use a circuit for SFE [Yao,GMW,BGW]
- Use k^2 private equality tests
 - Single inputs x,y ; return 1 iff $x = y$, 0 otherwise
 - $O(k)$ computation [NP]
- Diffie-Hellman based solutions [FHH99, EGS03]
 - Insecure against malicious adversaries
 - Depend on a “random oracle” assumption
- Our work: $O(k \ln \ln k)$ overhead.
 - “Semi-honest” adversaries – no RO assumption
 - “Malicious” adversaries – with RO assumption

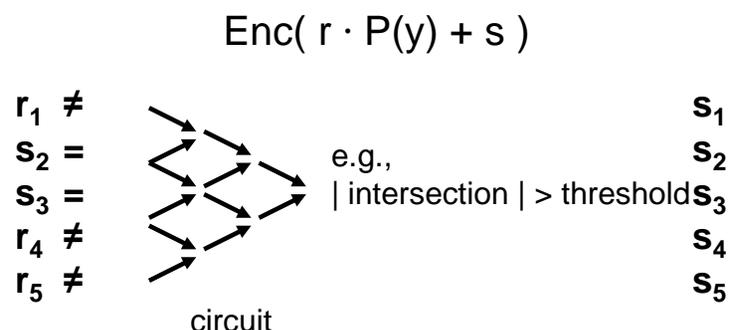
This talk...

- Overview
- Basic protocol in semi-honest model
- Efficient Improvements
- A little on...
 - Extending protocol to malicious model
 - Approximation bounds
 - Multi-party security
 - Fuzzy matching

Basic tool: Homomorphic Encryption

- Semantically-secure public-key encryption
- Given $\text{Enc}(M1)$, $\text{Enc}(M2)$, can compute
 - $\text{Enc}(M1+M2) = \text{Enc}(M1) \cdot \text{Enc}(M2)$
 - $\text{Enc}(c \cdot M1) = [\text{Enc}(M1)]^c$, for any constant cwithout knowing decryption key
- Examples: El Gamal, Paillier, DJ

Variant protocols...others



- $\forall y$, compute $r \cdot P(y) + s$, for $s \leftarrow \text{random}$
- Perform Yao circuit on decrypted values

Security (semi-honest case)

- Client's privacy
 - S only sees semantically-secure enc's
 - Learning about C's input = breaking enc's
- Server's privacy (proof via simulation)
 - Client can simulate her view in the protocol, given the output of $X \cap Y$ alone: she can compute the enc's of items in $X \cap Y$ and of random items.

Efficiency

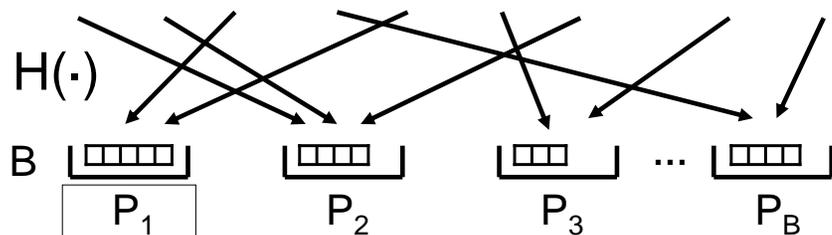
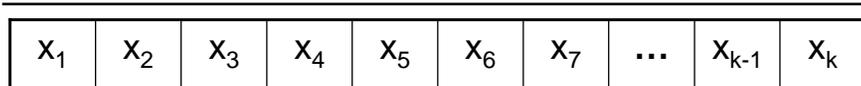
- Communication is $O(k)$
 - ✓ C sends k coefficients
 - ✓ S sends k evaluations on polynomial
- Computation
 - ✓ Client encrypts and decrypts k values
 - ✗ Server:
 - $\forall y \in Y$, computes $\text{Enc}(r \cdot P(y) + y)$, using k exponentiations
 - Total $O(k^2)$ exponentiations

Improving Efficiency (1)

- Inputs typically from a "small" domain of D values. Represented by $\log D$ bits (...20)
 - Use Horner's rule

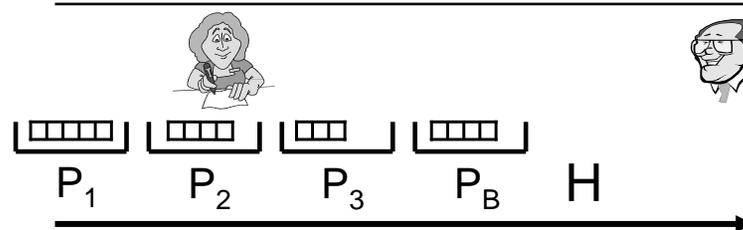
$$P(y) = a_0 + y (a_1 + \dots + y (a_{n-1} + y a_n) \dots)$$
 - That is, exponents are only $\log D$ bits
 - Overhead of exponentiation is linear in | exponent |
- Improve by factor of | modulus | / $\log D$
e.g., $1024 / 20 \approx 50$

Improving Efficiency (2): Hashing



- C uses PRF $H(\cdot)$ to hash inputs to B bins
- Let M bound max # of items in a bin
- Client defines B polynomials of deg M . Each poly encodes x 's mapped to its bin

Improving Efficiency (2): Hashing



$$\forall y, i \leftarrow H(y), r \leftarrow \text{rand}$$

$$\text{Enc}(r \cdot P_i(y) + y)$$

- Client sends B polynomials and H to server.
- For every y , S computes $H(y)$ and evaluates the single corresponding poly of degree M

Overhead with Hashing

- Communication: $B \cdot M$
- Server: $k \cdot M$ short exp's, k full exp's
 $(P_i(y))$ $(r \cdot P_i(y) + y)$
- How to make M small as possible?
 Balanced allocations [ABKU]:
 - H : Choose two bins, map to the emptier bin
 - $B = k / \ln \ln k \rightarrow M = O(\ln \ln k)$ ($M \leq 5$ [BM])
 - Communication: $O(k)$
 - Server: $k \ln \ln k$ short exp, k full exp in practice

This talk...

- Overview
- Basic protocol in semi-honest model
- Efficient Improvements
- A little on...
 - Extending protocol to malicious model
 - Approximation bounds
 - Multi-party security
 - Fuzzy matching

Malicious Adversaries

- Malicious clients
 - Without hashing: trivial. Parties use known a_0
 - With hashing
 - Verify that total # of roots (in all B poly's) is k
 - Solution using cut-and-choose
 - Exponentially small error probability
 - Still standard model
- Malicious servers
 - Privacy...easy:
S receives semantically-secure encryptions

Security against Malicious Server

- Correctness: Ensure that there is an input of k items corresponding to S's actions
- Problem: Server computes $r \cdot P(y) + y'$
- Solution: Server uses RO to commit to seed, then uses resulting randomness to "prove" correctness of encryption

Is Approximation easier?

- Represent inputs sets as k-bit vectors
0 0 1 1 1 0 0 1 0 1 0 1 0 0 1 0 1 0 1
- Approximate size of intersection (scalar product) with sublinear overhead? And securely?
- Lower bound:
 - Approximating $|X \cap Y|$ within $1 \pm \epsilon$ factor requires $\Omega(k)$ communication
 - True even for randomized algorithms
 - Proof: Reduction from Razborov's lower bound for Disjointness
- We provide secure approximation protocol

Multi-party intersection

- N parties: (N-1) clients, 1 leader
- $\forall y$, leader prepares (N-1) shares that XOR to y
- Each client performs intersection protocol with leader, learns random share of y
- Clients XOR (N-1) decrypted values
Recovers y iff $y \in |X_1 \cap X_2 \cap X_3 \cap \dots \cap X_N|$
- Nice communication flow

Fuzzy matching

- Databases are not always accurate or full
 - Errors, omissions, inconsistent spellings, etc.
- How to report a match iff entries similar?
 - Match in t out of T “attributes”
- Adaption of earlier protocol, but requires T choose t overhead

Open problems

- More computationally-efficient protocol?
- Malicious parties
 - Protocol secure in standard model?
 - Secure, efficient set cardinality protocol?
- Fuzzy matching
 - Efficient protocol needed?
 - Security in malicious model?