



Brice Minaud

email: brice.minaud@ens.fr





ENS



Techniques in Cryptography and Cryptanalysis Part 2: Oblivious Algorithms

"Techniques in Cryptography and Cryptanalysis": will cover (a choice of) important areas of cryptography.

- Zero-knowledge proofs
- Oblivious algorithms
- Lattices

Meta information

✓ done



2nd period



Big promises

Course contents:

- 1. Motivation (big promises).
- 2. Oblivious sorting.
- 3. Oblivious RAM.

Roadmap

Ad-hoc schemes General schemes

Outsourcing storage



Client

Examples:

- Company outsourcing customer/transaction info.
- Private messaging service.
- Trusted processor accessing RAM.

Server

Scenario: Client outsources storage of sensitive data to Server.

Outsourcing storage



Client

Adversary: honest-but-curious server. Security goal: privacy of data and queries.

Server = adversary

- **Scenario:** Client outsources storage of sensitive data to Server.

Privacy of data: just encrypt?

Symmetric Encryption.

Message space **M**, ciphertext space **C**, key space **K**.

Setup: Pick key $K \leftarrow_{\$} K$.

Encryption: encryption of $M \in \mathbf{M}$ is $C = Enc_{\mathcal{K}}(M) \in \mathbf{C}$.

Decryption: decryption of C is $M = \text{Dec}_{K}(C)$.

Perfect secrecy (sketch):

Let M_0 and M_1 be two arbitrary messages. **Perfect secrecy:** $Enc_{\mathcal{K}}(M_0) = Enc_{\mathcal{K}}(M_1)$. The equality is an equality of distributions. The randomness is over

the uniform choice of K.

IND-CCA: indistinguishability under Chosen-Ciphertext Attacks



7

Symmetric encryption: "proper" definition

Symmetric Encryption.

Message space **M**, ciphertext space **C**, key space **K**.

Setup: Pick key $K \leftarrow_{\$} K$.

Encryption: encryption of $M \in \mathbf{M}$ is $C = Enc_{\mathcal{K}}(M) \in \mathbf{C}$.

Decryption: decryption of C is $M = \text{Dec}_{K}(C)$.

Correctness: for all $M \in M$,

Dec_k(E

Security: IND-CCA (for instance

Caveats:

- nonces (\rightarrow mode of operation).

$$Enc_{\kappa}(M)) = M.$$

- Deterministic scheme cannot be IND-CCA. Need randomness, or

- "Security" above only covers **confidentiality**, not **integrity**.

Naive solution



Client

Naive solution: encrypt all data using symmetric encryption.

How can you fetch (or update) specific files?

Server

✓ Secure × Not functional



Naive solution, attempt #2



Client

Naive solution 2: encrypt each file individually.

What about security?

The server learns your access pattern.

Server

✓ Functional × Not secure

Some perspective: computing on encrypted data

<u>SNARKs</u>: prove arbitrary statements on encrypted data.





Is it okay to leak access pattern?

(No.)



Does leaking access pattern matter?





Imagine hospital storing patient information.

Connection with machine learning.

Example: range queries

- Sometimes searches for all patients with ages between a and b.
 - What can the server learn from the above leakage?









Foundational paper: Vapnik and Chervonenkis, 1971. Uniform convergence result.

Now a foundation of learning theory, especially PAC (probably approximately correct) learning.

Wide applicability.

Fairly easy to state/use. (You don't have to read the original article in Russian.)

VC Theory

Warm-up

Set X with probability distribution D. Let $C \subseteq X$. Call it a concept.



$Pr(C) \approx \frac{\#points in C}{\#points total}$

Sample complexity: to measure Pr(C) within ε , you need O(1/ ϵ^2) samples.

Approximating a Concept Set

Now: set \mathcal{C} of concepts. Goal: approximate their probabilities simultaneously.



The set of samples drawn from *X* is an **ε-sample** iff for all C in C:

$$\left| \Pr(C) - \frac{\# \text{points in } C}{\# \text{points total}} \right| \le \epsilon$$

ε-sample Theorem



How many samples do we need to get an ε -sample whp?

Union bound: yields a sample complexity that depends on $|\mathcal{C}|$.

V & C 1971:

If C has VC dimension d,

then the number of points to get an *\varepsilon*-sample whp is

$$O(\frac{d}{\epsilon^2}\log\frac{d}{\epsilon}).$$

Does not depend on $|\mathcal{C}|$!

VC Dimension

Remaining Q: what is the VC dimension?

A set of points is **shattered** by *C* iff:



- every subset of S is equal to $C \cap S$ for some C in \mathcal{C} .



VC dimension of \mathcal{C} = largest cardinality of a set of points in X that is shattered by \mathcal{C} .

E.g. VC dimension of ranges is 2.

What typically matters is just that VC dim is finite.

VC Dimension



Statistical Reconstruction





VS.



Direct statistical attack



 \rightarrow gives estimate of probability it's hit by uniform query. **2.** deduce estimate of its value by "inverting" **f**.

Direct statistical attack



This is an ε -sample!

X = ranges $\mathcal{C} = \{\{ranges \ni x\}: x \in [1,N]\}$

so we need O($\epsilon^{-2} \log \epsilon^{-1}$) queries.

After O(ϵ^{-4} log ϵ^{-1}) queries, the value of all records is recovered within εN .

- **Step 1**: for all records, estimate prob of the record being hit.
- **Step 2**: because **f** is quadratic, "inverting" **f** adds a square.

Order Reconstruction







Problem Statement



distribution.

- What can the server learn from the above leakage?
- This time we don't assume i.i.d. queries, or knowledge of their

Range Query Leakage

Query A matches records a, b, c. Query B matches records b, c, d.



Then this is the only configuration (up to symmetry)!

 \rightarrow we learn that records b, c are between a and d.

We learn something about the order of records.

Range Query Leakage

Query A matches records a, b, c. Query B matches records b, c, d. Query C matches records c, d.



Challenges:

- as more queries are observed?

How do we extract order information? (What algorithm?) How do we quantify and analyze how fast order is learned

Short answer: there is already an algorithm!

Long answer: PQ-trees.

X: linearly ordered set. Order is unknown.

Can be updated in linear time.

- Challenge 1: the Algorithm

- You are given a set S containing some intervals in X.
 - A **PQ tree** is a compact (linear in |X|) representation of the set of all permutations of X that are compatible with S.

Challenge 2a: quantify order learning

exact value of every record.

all values within εN .

 $\varepsilon = 0.05$ is recovery within 5%. $\varepsilon = 1/N$ is full recovery.

- **Strongest goal: full** database reconstruction = recovering the
- **More general:** approximate database reconstruction = recovering

("Sacrificial" recovery: values very close to 1 and N are excluded.)

Challenge 2b: analyze query complexity



Intuition: if no query has an endpoint between a and b, then a and b can't be separated.

 $\rightarrow \varepsilon$ -approximate reconstruction is impossible.

You want a query endpoint to hit every interval $\geq \varepsilon N$. Conversely with some other conditions it's enough.

Heavy sweeping of details under rug.



The set of samples drawn from X is an ε -net iff for all C in \mathcal{C} :

 \rightarrow Number of points to get an ε -net whp:

VC Theory saves the day (again)

ε-samples: the ratio of points hitting each concept is close to its probability.

What we want now: if a concept has high enough probability, it is hit by at least one point.

 $Pr(C) \ge \epsilon \Rightarrow C$ contains a sample

$$O\left(\frac{d}{\epsilon}\log\frac{d}{\epsilon}\right)$$

Access pattern leakage: conclusion

Say patient age has N possible values (e.g. N = 100)...

Full order reconstruction: O(*N* log *N*) queries. (NB: this is optimal.)

Age data: can infer value from order (if all ages are present)... In this setting, encryption was ultimately useless.

Very rough summary :

- **Approximate order reconstruction** (within εN): O($\varepsilon^{-1} \log \varepsilon^{-1}$) queries!

- highly structured queries
 - \Rightarrow low VC dimension
- \Rightarrow learn data with few queries



	-
	_
_	

Other examples

If adversary can observe queries to tables, AES is **broken**. tables, also broken.

Two issues:

- Leaking access pattern can be (very) damaging.
- Many settings leak access pattern, completely or partially.

hypervisors), etc. See also: side-channel attacks.

Suppose you implement AES using lookup tables (for S-boxes).

If adversary can observe cache misses from access to AES S-box

Cloud storage, trusted enclaves, cache attacks (incl.

Oblivious algorithms




Naive solution, attempt #2



Client

Naive solution 2: encrypt each file individually.

The server learns the client's access pattern.

Remark: the "client" could be an algorithm.

Server

✓ Functional × Not secure





Server stores *N* files.

Client fetches file *i*.

Security: Server learns *nothing* about *i*.

Efficiency: ORAM algorithm only queries O(log N) files.

Bonus feature: server performs no computation. Acts like a RAM.

Magic Claim

Oblivious algorithm: definition

Oblivious algorithm: an algorithm *A* is oblivious iff for any two inputs *x* and *y*, the memory accesses of *A* on input *x*, and *A* on input *y*, are indistinguishable.

Oblivious Sorting





Sorting algorithms

Oblivious algorithm: an algorithm *A* is oblivious iff for any two inputs *x* and *y*, the memory accesses of *A* on input *x*, and *A* on input *y*, are indistinguishable.

Which of the following algorithms are oblivious? (assuming inputs are arrays of fixed size.)

1. Bubble Sort.

2. Quick Sort.

3. Merge Sort.



🗙 no

🗙 no

Sorting obliviously

Basic operation: sorting two elements.

output (y,x).



Compare and swap: on input (x,y), if x < y, output (x,y), else

Bubble Sort



Sorting network

Comparator network: A comparator network is an algorithm that consists in a sequence of compare-and-swaps ("comparators") between fixed inputs.

Can be represented in this form:



Sorting network: A sorting network is a comparator network that correctly sorts its input (for all possible inputs).

Remark: testing whether a comparator network is a sorting network is co-NP-complete.

Size and depth

The size of a comparator network is its number of comparators.

The depth (or "critical path") of a comparator network is the maximum number of comparators that an input value can go through.

It is also the number of steps in a parallel computation of the network.



Example: comparator network of size 4 and depth 3.

→ Bubble sort is a sorting network of size $O(n^2)$ and depth O(n).

Can we do better?

Proposition: A sorting network must have size $\Omega(n \log n)$.

Proof. A network with *k* comparators can permute its input sequence in at most 2^k different ways. For a sorting network, we must have $2^k \ge n!$ By Stirling's formula, this yields $k = \Omega(n \log n)$.

Bitonic sort: size $O(n \log^2 n)$. Most efficient in practice.

The 0-1 principle

0-1 Principle: a comparator network (on *n* inputs) is a sorting network iff it correctly sorts all 2ⁿ possible binary inputs.

Proof of the claim: induction on comparators.

but they are in the opposite order in the output.

Define f(x) = 0 if $x \le x_i$, 1 otherwise. We have $f(x_i) = 0$ and $f(x_i) = 1$, but their order is reversed by the network when inputting $(f(x_1), \ldots, f(x_n))$. Hence the network does not correctly sort all binary sequences.

- *Proof.* Let f be a non-decreasing function: $x \le y$ implies $f(x) \le f(y)$.
- *Claim:* If a comparator network has input (x_1, \ldots, x_n) and outputs
- (y_1,\ldots,y_n) , then on input $(f(x_1),\ldots,f(x_n))$, it must output $(f(y_1),\ldots,f(y_n))$.

Now assume a comparator network is not a sorting network. Then there exist an input (x_1, \ldots, x_n) and some indices *i*, *j*, such that $x_i < x_j$

Bitonic sequences

Bitonic sequence: A sequence of values is **bitonic** iff:

- It is increasing, then decreasing.
- Or it is a circular shift of the previous case.

Example: bitonic sequences of 0 and 1's are those of the form $0^{a}1^{b}0^{c}$ and $1^{a}0^{b}1^{c}$.

number of inputs *n*, composed of comparators

Half-cleaner for n = 8:



Half-cleaner

Half-cleaner: A half-cleaner is a comparator network for an even

(1,*n*/2+1), (2,*n*/2+2), ..., (*n*/2,*n*)

Half-cleaner

must be all 0's or all 1's. That half is called clean.



Key property of a half-cleaner: if the input is bitonic, then both halves of the output are bitonic. Moreover, one of the two halves

Bitonic sorter



The bitonic sorter correctly sorts all **bitonic** inputs.



Batcher's sort



Batcher's sort correctly sorts all binary inputs, hence all inputs.



A half-cleaner has size H(n) = n/2. The bitonic sorter has size $S(n) = H(n) + 2S(n/2) = O(n \log n)$. Batcher's sort has size $B(n) = S(n) + 2B(n/2) = O(n \log^2 n)$.

model, only need $O(\log^2 n)$ steps.

 \rightarrow Sorting algorithms used in GPUs.

Ajtai, Komlós, Szemerédi (STOC '83): there exists a sorting network of size O(n log n).

Unfortunately, completely impractical.

Efficiency

- The depth of Batcher's sort is O(log² n): in a parallel computation

Oblivious RAM







Generalizing

So far...

Traditional efficient sorting algorithms were not oblivious.

 \rightarrow created new efficient oblivious sorting algorithm.

This is what Oblivious RAM (ORAM) does.

Disclaimer: does not hide number of accesses.

- Can we do this generically?
- Take any algorithm \rightarrow create oblivious version, with low overhead.

Reminder: Oblivious RAM



Client wants to do queries $q_1, q_2, ..., q_n$.

Each q_i is either:

- read(a): read data block at address a;
- write(a,d): write data block d at address a.

Reminder: Oblivious RAM



ORAM algorithm C (or ORAM "compiler"): transforms each query q by the client into one or several read/write queries C(q) to server.

Correctness: C's response is the correct answer to query q.

 $(C(r_1),\ldots,C(r_k))$ are indistinguishable.

Obliviousness: for any two sequences of queries $q = (q_1, \ldots, q_k)$ and $r = (r_1, ..., r_k)$ of the same length, $C(q) = (C(q_1), ..., C(q_k))$ and C(r) =

Trivial ORAM



Trivial ORAM: read and re-encrypt every item in server memory.

Security: trivial.

Efficiency: every client query costs O(n) real accesses \rightarrow overhead is O(n).

A non-trivial ORAM must have:

- Client storage o(n).
- Query overhead o(n).

Suppose client wants to do queries q_1, q_2, \ldots, q_n . Each q_i is to read or write a block of memory.

Assume the client does not store any memory block.

memory.



- Some observations
- For each *q_i*, the ORAM has to do some access(es) to the server

 q_1 and q_2 must access at least 1 data

qn ...



(q_1,q_2) and (q_3,q_4) must access at least 2 data blocks in common.



etc...

. . .

accessed (i < j).



For each memory access done by q_i , let's "assign" that access to the node $q_i \wedge q_j$, where q_i is the last time the same address was

> At least this many accesses are assigned to this node.

accessed (i < j).



For each memory access done by q_i , let's "assign" that access to the node $q_i \wedge q_j$, where q_i is the last time the same address was

\rightarrow Memory accesses in every node are now unique to that node.

(Say $n = 2^k$ for some k.)



A lower bound

Goldreich & Ostrovsky '96 (again):

G&O's proof under assumptions:

- Client memory O(1).
- Statistically secure ORAM.
- "Balls and bins" model.

What we just saw: stronger proof by Larsen & Nielsen '18. (Computational security, less restrictive cell probe model, online).

Secure ORAM must have overhead $\Omega(\log n)$.

Proof sketch (Goldreich-Ostrovsky proof)

Each item i = colored ball. At start:



Suppose client wants to make queries for balls b_1, \ldots, b_q .

put ball, take ball, nothing.

all n^q possible query sequences (b_i).

$$\Rightarrow O($$

$$\Rightarrow f(q) =$$

- Server: n balls + extra room
- \rightarrow ORAM makes accesses $a_1, \ldots, a_{f(q)}$. (Includes Setup accesses.)
- Each server access, ORAM can do O(c) operations: exchange ball,
- Statistical security \rightarrow access sequence (a) must be compatible with
- But only O(c)^{f(q)} possible sequences of balls held by client, hence O(c)^{f(q)} query sequences compatible with given access sequence.
 - $(C)^{f(q)} \ge n^q$
 - *q* · Ω(log *n*)

Query overhead: how many queries to the server are made in C(q)for each client query q, amortized (= on average).

Family of constructions

- 1. Square-root ORAM
- 2. Hierarchical ORAM
- 3. Tree ORAM

 $polylog(x) = poly(log(x)) = O(log^{c}(x))$ for some constant c.

Other efficiency metrics: client memory size, number of roundtrips in C(q), time complexity of C...

Roadmap

Here, $n = \max$ memory size = max number of items (address, data).

<u>Overhead</u>	<u>Feature</u>
Õ(n ^{1/2})	Simple
O(polylog <i>n</i>)	Best in theory
O(polylog <i>n</i>)	Best in practice

Square-root ORAM



- The **main memory** stores *n*+*s* items:
- Real items: with addresses in [1,*n*], real data.
- **Dummy items:** with addresses in [n+1,n+s], random data.
- For now, stash contains s items with all-zero address and data.





main memory

1. Client chooses permutation π over [1,*n*+s]. Item *i* will be stored at location $\pi(i)$ in the main memory.

2. Client encrypts everything, and sends to server.

Server view:

n+s encrypted items

Remark: we are assimilating client with ORAM algorithm.

Setup



s enc. items



- **3.** Rewrite whole **stash** to server.

Lookup

2. Read/rewrite **location** $\pi(n+t)$ in main memory. $t \leftarrow t+1$.



items (t > s).

Can only happen after s iterations.

- Solution: after s iterations of lookup, perform refresh: • Client chooses new permutation π '.
 - Moves item *i* to location $\pi'(i)$ in main memory.
 - Empties stash.
- \rightarrow equivalent to fresh setup with π '.
- \rightarrow can do s iterations again...

How do you move item *i* to location $\pi'(i)$ obliviously?

Refresh

Lookup can fail in two ways: stash is full, or run out of fresh dummy

Oblivious sorting!

Refresh via oblivious sort

Server memory after s lookups...

main memory

n real items (some outdated), s

n real items (some duplicate

n real items + some empty

n+s items sorted with

main memory

		stash		
dummies real items,		empty items		
	Oblivious sort with π ⁻¹			
es)	s dummies		empty items	
	Erase outdated duplicates			
y	s dummies		empty items	
	Oblivious sort with π'			
π'		empty items		
		stash 72		


Setup: server sees:

Lookup: server sees:



Refresh: server sees:

1 oblivious sort, 1 linear scan, 1 oblivious sort.

Remark: computationally secure. Essentially statistically secure, except for encryption, and pseudo-random permutation π .

Security

73



Overhead. Lookup costs O(s). **Refresh** costs O(*n* polylog *n*), happens every s lookups. **Total overhead** (amortized): Setting $s = n^{1/2} \log n$, and using Batcher sort:

Server memory: O(n).

Client memory: O(1).

operations.

Remark: memory measured in number of items. Item size assumed to be $\Omega(n)$ bits, which is also $\Omega(\lambda)$ if $n \ge \lambda$.

Efficiency

- $O(s + n/s \cdot polylog(n))$

 - $O(n^{1/2} \log n)$

Need encryption key + key for pseudo-random π + few items during

Hierarchical ORAM





Hierarchical ORAM

Goldreich and Ostrovsky '96:

- Square-root ORAM, overhead $\tilde{O}(n^{1/2})$.
- Secure ORAM must have overhead $\Omega(\log n)$.

But also: hierarchical ORAM, overhead O(log³ n). \rightarrow Spawned whole construction family of ORAMs.

Interesting because:

- First ORAM with polylog overhead.
- Basis for the recent construction of optimal ORAM with overhead O(log n).

Open problem for 20+ years, solved by Asharov et al. '18, based on Patel et al. '18.

Hash function H: $\{0,1\}^* \rightarrow [1,n]$. Want to store *n* items into *n* buckets according to H.

items hash H buckets



Proof: Probability that given bucket receives more than k items is $exp(-\Omega(k^2))$ by Chernoff bound. Union bound over all buckets:

 $n \cdot \exp(-C \cdot \log^2 n)) = n^{1 - C \cdot \log n} = \operatorname{negl}(n).$

Hashing



Buckets of size log *n* suffice for negligible probability of overflow.

Oblivious hashing

Want to do the assignment obliviously... Suppose we have items + empty buckets all in server memory.



- *Sketch:* **1**. obliviously sort items according to H.
 - 2. Put each item into own bucket.
 - 3. Scan all buckets, pushing content of each bucket into next bucket if next bucket has same hash value.
 - 4. Obliviously sort *buckets* to delete empty buckets.

Assignment can be done obliviously in *n* log² *n* operations.

78



Server memory arranged into log n levels. Each level k is an (oblivious) hash table for 2^k items.

At start:

All items are in last level. Other levels contain dummies.



To access item item *i*:

- 1. Access each level k at location $H_k(i)$ until item is found. 2. Access remaining levels at uniformly random location. 3. Insert item at level 1. (Potentially with new value.)



Remark: whenever accessing level 1, entire level is read + rewritten.



To maintain invariant that level k stores $\leq 2^k$ items:

into level k+1, using fresh hash function.

Invariant is preserved:

Level k receives at most 2^{k-1} items every 2^{k-1} lookups. And empties its content every 2^k lookups.

are no duplicate items in the same level.

Reshuffling

- Every 2^k lookups, the (non-dummy) items of level k are shuffled
- If an item appears twice, newest version (from earliest level) is kept.

- *Remark:* last level is never full, because it can hold *n* items, and there

Setup: server sees log *n* hash tables:



Lookup: server sees:



Key fact: no item is ever read twice from the same level with the same hash function.

Security



Overhead.

- Level k is reshuffled every 2^k lookups.
- Each reshuffle costs: O($2^k \log^2 n$).
- \rightarrow Amortized cost for level k: O(log² n).
- \rightarrow Total amortized cost of reshuffles: O(log³ n).
- \rightarrow Total amortized overhead: O(log³ n) + O(log² n) = O(log³ n).
- Server memory: O(n log n).
- Client memory: O(1).

Server memory can be reduced to O(n) using cuckoo hashing.

Efficiency

Cuckoo hashing

Initial design mainly motivated by real-time systems...

Idea:

items

hashes H₁, H₂

m=O(n) cells



Each item *i* can go into one of two cells $H_1(i)$ or $H_2(i)$.

- "Bucket" hashing had total storage O(n log n), and lookup O(log n). Cuckoo hashing has storage $(2+\epsilon)n = O(n)$, and lookup 2 = O(1).

The cuckoo graph

Picture graph with cells = nodes, item $i = \text{edge } H_1(i) - H_2(i)$.



The cuckoo graph

Picture graph with cells = nodes, item $i = edge H_1(i) - H_2(i)$. Orient edge towards where item is stored.



its other possible cell. Repeat until unoccupied cell is reached.

To insert item *i*: try cell $H_1(i)$. If occupied, move occupying item into

The cuckoo graph

Picture graph with cells = nodes, item $i = edge H_1(i) - H_2(i)$. Orient edge towards where item is stored.



its other possible cell. Repeat until unoccupied cell is reached.

To insert item *i*: try cell $H_1(i)$. If occupied, move occupying item into



has at most one cycle.

Moreover, with *n* edges and $m = (2+\varepsilon)n$ nodes...

- The previous fact holds with high probability.
- Expected size of a connected component is O(1).

 \rightarrow Expected insertion time is O(1)!

Remark: Probability of failure can be made negligible by adding a stash.



- **Theorem:** assignment is possible iff every connected component

Tree ORAM







Hierarchical ORAM family leads to recent optimal construction. But huge constants. Never used in practice.

What is actually used:

Overhead: $O(\log^3 n)$. Worst-case (no need to amortize).

In practice: easy to implement, efficient.

We will see Simple ORAM, member of the Tree ORAM family.

Tree ORAM

by Shi et al. '11 Tree ORAM



Server-side memory is a full binary tree with $log(n/\alpha)$ levels. Each node contains log *n* blocks. Each block contains $\alpha = O(1)$ (possibly dummy) items.



At start:

Each block b is stored in a uniformly random leaf Pos(b). "Position map" Pos() is stored on the client.

- Items are grouped into blocks of α items, item *i* into block $b = \lfloor i/\alpha \rfloor$.
- 92 **Invariant:** block b will always be stored on the branch to Pos(b).



To access item *i* from block *b*: 2. Update Pos(b) to new uniform leaf. 3. Insert *b* at root. (Possibly with new value.)

- 1. Read every node along branch to Pos(b). Remove b when found.



After every lookup

1. Pick branch to uniformly random leaf. that block b must remain on branch to Pos(b)).

2. Push every block in the branch as far down as possible (preserving

Security

Setup: server sees full binary tree of height log (n/α) . Each node is encrypted, same size.



Full read/rewrite along 2 branches to uniformly random leaves.

Works as long as no node overflows. **Setup,** no overflow: same argument as bucket hashing.

Lookup + eviction, no overflow (sketch): Let K be the number of blocks per node (we had $K = \log n$). Pick arbitrary node x at level L. least K.

below that child c must be at least K/2.

Both events have the same probability (namely 2^{-L}).

Remark: we cheat a little by setting $K = \log n$.

- Why does that work? (sketch)
- For x to overflow, number of blocks whose Pos is below x must be at
- \rightarrow For one of the two children of x, number of blocks whose Pos is
- → This implies event [Pos of new block is below c] happens K/2 times, without event event eviction branch includes c happening at all.
- Deduce overflow probability is $\leq 2^{-K/2}$. Negligible for $K = \omega(\log n)$.

Efficiency of basic construction

Overhead. Each lookup, read two branches, total $O(\log^2 n)$ items. Server memory: $O(n \log n)$. Client memory: $O(n/\alpha)$. (oops)

The position map

Still a large gain, if item size is much larger than $\log(n/\alpha)$ bits.

To reduce client memory:

"Recursive" construction:

Client needs new position map for server-side position map... **Key fact:** it is α times smaller!

Repeat this recursively $\log_{\alpha}(n)$ times. In the end:

- Client position map becomes size O(1).
- Server stores $\log_{\alpha}(n)$ position maps, each $\alpha \times$ smaller than last.
- Each lookup, $\log_{\alpha}(n)$ roundtrips to query each position map.

The client stores position Pos: $[1,n/\alpha] \rightarrow [1,n/\alpha]$, size $n/\alpha = \Theta(n)$.

- Store position map on server. Obliviously!

Efficiency of recursive construction

Overhead.

Each lookup, O(log n) recursive calls, ecah of size O(log² n).

 \rightarrow O(log³ *n*) overhead.

Server memory: O(n log n).

Client memory: O(1).

Original Tree ORAM had more complex eviction strategy and analysis, better efficiency.

Path ORAM:

- Client has a small stash of blocks.
- Blocks are evicted along the same branch as item was read.
- Can use nodes as small as K = 4 blocks!

Variants

Searchable Encryption







Encrypted search:

- Client stores encrypted database on server.
- Client can perform search queries.
- Privacy of data and queries is retained.
- Example: private email storage.

Dynamic SSE: also allows update queries.

Searchable Symmetric Encryption

Two databases:

Document database:

Encrypted documents d_i for $i \leq D$.

(Reverse) index database DB:

that d_i contains W.

- Pairs (W,i) for each keyword W and each document index i such

 $\mathsf{DB} = \{(\mathbf{W}, \mathbf{i}) : \mathbf{W} \in \mathbf{d}_{\mathbf{i}}\}$

A simple solution

Put everything into ORAM.

Secure.

(Up to leaking lengths of anwers.)

Inefficient.

(In certain cases, such as Enron email dataset or English Wikipedia, some studies suggests trivial ORAM would be most efficient.)

Idea of Searchable Encryption: allow some leakage. Gain in efficiency.

How to capture leakage

only the number of accesses.

• Formally: secure iff there exists a simulator, which on input number of accesses, outputs a set of accesses indistinguishable from real algorithm.

simulator knowing only the output of a leakage function L.

• Formally: secure iff there exists a simulator, which on input the output of the leakage function, outputs a set of accesses indistinguishable from real algorithm.

ORAM security: accesses can be simulated by a simulator knowing

Searchable encryption security: accesses can be simulated by a

(Leakage function takes as input the database and all operations.)

Security Model

