

Free semi-group ciphers

Alexander Towell
atowell@siue.edu

Abstract

We define the semantics of abstract data type for the *oblivious set*. We demonstrate a theoretical data structure, denoted the *Singular Hash Set*, that provides an optimal implementation of the oblivious set with respect to *entropy* and *space complexity*. Finally, we show how to use the *Bloom filter* and *Perfect Hash Filter* to implement oblivious sets and compare them to each other and the theoretically optimal *Singular Hash Set*.

Contents

1	Introduction	1
2	Cipher sets	2
2.1	Notation	3
2.2	Oblivious object types	3
3	Cipher Boolean algebras	4
4	Cipher Boolean ring	6
5	The <i>Singular Hash Set</i>	7
5.1	Time and space complexity	9
5.2	Probability distributions	11
5.3	Relaxing optimality	13
5.4	Boolean algebra	16

1 Introduction

The *oblivious set*[?] is a fundamental data structure that may be used to construct other oblivious object types, like secure indices for Boolean Encrypted Search[?]

An *oblivious set* is an oblivious object type over $(2^{\mathcal{U}}, \in)$, which means that the oblivious set represents values in the set $2^{\mathcal{U}}$ over the binary *member-of* predicate $\in: \mathcal{U} \times 2^{\mathcal{U}} \mapsto \{0, 1\}$.

Definition 1.1 (n -fold Cartesian product). *Let X_1, \dots, X_n denote n sets. The set*

$$X_1 \times \dots \times X_n := \{(x_1, \dots, x_n): x_1 \in X_1 \wedge \dots \wedge x_n \in X_n\} \quad (1)$$

is called the Cartesian product of sets X_1, \dots, X_n .

A shorthand notation for the Cartesian product $X \times X \times X$ is denoted by X^3 .

Suppose we have sets X_1, X_2, \dots, X_n and a surjective pairing function

$$f: \mathbb{N}^2 \mapsto \mathbb{N}. \quad (2)$$

Recursively, the n -tuple encoding function $g_n: \mathbb{N}^n \mapsto \mathbb{N}$ may be defined as

$$g_2(x_1, x_2) := f(x_1, x_2) \quad (3)$$

$$g_n(x_1, \dots, x_n) := f(x_1, g_{n-1}(x_2, \dots, x_n)). \quad (4)$$

Assuming we have a serializer $h_j: X_j \mapsto \mathbb{N}$ for $j = 1, \dots, n$, an encoder for tuples of type $X_1 \times \dots \times X_n$ is given by

$$\text{encode}(X_1, \dots, X_n) := g_n(h_1(X_1), \dots, h_n(X_n)). \quad (5)$$

Now, we may use any random approximate set over the natural numbers to represent any n -tuple relation.

A *regular function* over some abstract data type X behaves the same way if given any data structure that implements the behavior of X .

1. Approximate membership tests of specific elements may be performed with a false positive rate ε and a false negative rate ω . Note that there is no way to efficiently iterate over the elements.
2. The cardinality may be estimated to be within some range. The the degree of uncertainty can be made arbitrarily large entropy at the expense of its space complexity.
3. Set-theoretic operations like union, intersection, and complement generate oblivious sets that approximate the true operation as a function of the false positive and false negative rates and the degree of similarity between the exact sets under consideration.

In section 2, we precisely define the oblivious set.

In ??, we derive the probabilistic model of *oblivious sets*. In ??, we derive the *entropy* of *oblivious sets*. In ??, we derive estimators of properties of *oblivious sets*. In section 5, we provide a theoretically optimal implementation of the oblivious set.

2 Cipher sets

A set is given by the following definition.

Definition 2.1. *A set is an unordered collection of distinct elements from a universe of elements.*

A countable set is a *finite set* or a *countably infinite set*. A *finite set* has a finite number of elements. For example,

$$S = \{1, 3, 5\}$$

is a finite set with three elements. A *countably infinite set* can be put in one-to-one correspondence with the set of natural numbers

$$\mathbb{N} = \{1, 2, 3, 4, 5, \dots\}. \quad (6)$$

The cardinality of a set S is a measure of the number of elements in the set, denoted by

$$|S|. \quad (7)$$

The cardinality of a *finite set* is a non-negative integer and counts the number of elements in the set, e.g.,

$$|\{1, 3, 5\}| = 3.$$

Informally, a cipher set of \mathcal{S} , denoted by $\check{\mathcal{S}}$, provides a *confidential* in-place binary representation such that very little information about \mathcal{S} is disclosed.

Note that an oblivious set where elements are from the universe \mathcal{U} permits membership tests on elements in \mathcal{U} . These elements, and the universe \mathcal{U} , are *not* necessarily oblivious types. It is also possible to have an oblivious set over oblivious elements, where the elements have their own set of separate constraints, e.g., an oblivious type X in which only the less-than and equality predicates are possible. The oblivious set, of course, only allows *membership* tests to be performed on elements over $\check{\mathcal{U}}$.

2.1 Notation

The binary set $\{0, 1\}$ is denoted by \mathcal{B} . The set of all bit (binary) strings of length n is therefore \mathcal{B}^n . The cardinality of \mathcal{B}^n is

$$|\mathcal{B}^n| = 2^n. \quad (8)$$

The set of all bit strings of length n or less is denoted by $\mathcal{B}^{\leq n}$, which has a cardinality $2^{n+1} - 1$. The countably infinite set of all bit strings, $\mathcal{B}^{\leq \infty}$, is also denoted by \mathcal{B}^* .

The bit length of an object x is denoted by

$$\ell(x), \quad (9)$$

e.g., the bit length of any $x \in \mathcal{B}^n$ is $\ell(x) = n$.

A (convenient) one-to-one correspondence between \mathcal{B} and \mathbb{N} is given by the following definition.

Definition 2.2. Let the set of bit strings \mathcal{B}^* and the set of natural numbers \mathbb{N} have the bijection given by

$$(b_1 b_2 \cdots b_m) \longleftrightarrow 2^m + \sum_{j=1}^m 2^{m-j} b_j. \quad (10)$$

We denote the mapping of a bit string (or natural number) x by x' .

An important observation of this mapping is that a natural number n maps to a bit string n' of length $\ell(n') = \lfloor \log_2 n \rfloor$.

2.2 Oblivious object types

A *type* is a set and the elements of the set are called the *values* of the type. An *abstract data type* is a type and a set of operations on values of the type. For example, the *integer* abstract data type is defined by the set of integers and standard operations like addition and subtraction. A *data structure* is a particular way of organizing data and may *model* one or more abstract data types.

Suppose we have an abstract data type denoted by T with a set of operators \mathcal{F} denoted the *computational basis* that are (at least partially) functions of T . We denote that an *object* x in computer memory models T by $T(x)$.

An *oblivious* object type[?] that models T is a related type denoted by \check{T} that provides guarantees about what can be learned about an object \check{x} by observing its representation, memory access patterns, or any other possible correlations.

Optimally, the only information that can be learned about \check{x} is given by the well-defined *behavior* of the the abstract data type T on a subset of its computational basis \mathcal{F} .¹ For instance, say an operator $g: T \mapsto \mathcal{B}$ is not in the computational basis of \check{T} , i.e., $g: \check{T} \mapsto \mathcal{B}$ is undefined. Suppose $g(x) = 1$ for a particular object of T and half of the inputs map to 1 (and half map to 0). If the only information we have about x is given by \check{x} , then $P[g(\check{x}) = 1] = 0.5$, i.e., we can do no better than a random guess.

We consider random approximate cipher sets of the type $\boxed{2^{\check{X}}}^{\pm}$ where X is the underlying universal set of elements, i.e., cipher sets whose bit string representations are uncorrelated with its specific members over cipher elements, frequently denoted *trapdoors*. This is opposed to, say, $2^{\check{X}}$ in which we have a regular set over ciphers elements of type X , or $\boxed{2^{\check{X}}}$ cipher sets over elements of type X .

It is generally easier to understand what we mean by this with respect to a *computational basis*, e.g., the *singular hash set* has a computational basis given by

$$\in: \boxed{2^{\check{X}}} \times \check{X} \xrightarrow[\square]{+} \mathcal{B}. \quad (11)$$

Other operations, like set-theoretic operations, may be composed using this computational basis. Such a composition may contain information about the objective sets, but an explicit composition of cipher sets is still generally considered to be a cipher set.

A more confidential variation of a cipher set is given by

$$\in: \boxed{2^{\check{X}}} \times \check{X} \xrightarrow[\square]{+} \check{\mathcal{B}}, \quad (12)$$

i.e., its membership predicate returns a Boolean cipher. Since *representational equality* implies *equality*, there is always, at least implicitly, a predicate for the identity relation of the type

$$=: \boxed{2^{\check{X}}} \times \boxed{2^{\check{X}}} \mapsto \mathcal{B}. \quad (13)$$

However, equality does not necessarily imply representational equality, and so there is also an opportunity for a predicate of the type

$$=: \boxed{2^{\check{X}}} \times \boxed{2^{\check{X}}} \xrightarrow[\square]{+} \check{\mathcal{B}}, \quad (14)$$

but this may only be practical for relatively small cipher sets. If the cipher sets are random approximate sets over a *countably infinite* universe, then *equality* implies *representational equality*.

3 Cipher Boolean algebras

Cipher value types may come in many shapes. They may be parametric types, e.g., \check{T} is a cipher type of plain type T and they may provide differing levels of *confidentiality*, whether as a simple substitution cipher, homophonic cipher, or homomorphic encryption.

In what follows, we consider several ciphers sets, one in which is a Boolean algebra that provides for all set relations like *membership*, *identity*, and *subset*, and the other is a group over symmetric difference and intersection but only offers the *identity* relation (or only one its representations in a, say, homophonic encryption system). Both, however, model their relations as predicate functions which return *plaintext* Booleans rather than cipher variations, i.e., \mathcal{B} as opposed to the more confidential $\check{\mathcal{B}}$.

NOTE: If using something like homophonic encryption, then identity is no longer maintained. A cipher set $\check{\mathcal{A}}$ sent in one round may be different is *incomparable* to a cipher set $\check{\mathcal{A}}$ sent in another round.

¹Frequently, some operations may reveal too much information about the values and so are excluded from the computational basis of the oblivious type.

NOTE 2: When, say, including bigrams (for whatever reason), if we want to be able to construct a cipher bigram from two cipher elements on the untrusted system, then in the cipher map or set on the untrusted system, all of these combinations must be included separately (unless rate-distorted, in which case errors). Since this may limit its value, we can not support these operations and just go

and we wish with these cipher sets, then identity is no longer provided for by the information in the sets alone. See *approximate maps* or *approximate sets*.

Say we have the Boolean algebra

$$A := \left(2^{\mathcal{B}^*}, \cup, \cap, \overline{}, \emptyset, \mathcal{B}^* \right) \quad (15)$$

where \mathcal{B}^* is the free semigroup over \mathcal{B} , which is the closure over concatenation $+$: $\mathcal{B}^* \times \mathcal{B}^* \mapsto \mathcal{B}^*$.

Since \mathcal{B}^* can encode any tuple of elements, and $2^{\mathcal{B}^*}$ can encode any sum type of tuples, this Boolean algebra may be a reasonable model for hashing algebraic data types, while still permitting some operations on the resultant hash.

Consider the Boolean algebra

$$B := (\mathcal{B}^m, \vee, \wedge, \neg, 0^m, 1^m) \quad (16)$$

and suppose we wish to provide a homomorphism of type $A \mapsto B$ that, approximately, preserves relations while *obscuring* or *encrypting* the identities.

A *cryptographic hash function* is given by the following definition.

Definition 3.1. The hash function $h: \mathcal{B}^* \mapsto \mathcal{B}^m$ maps (hashes) bit strings of arbitrary-length to bit strings of fixed-length m and approximates a random oracle over its codomain.

Definition 3.2. The homomorphism $F: A \mapsto B$ is defined as

$$F(\beta) := \begin{cases} h(x) & \beta \in \mathcal{B}^* \\ \vee & \beta = \cup \\ \wedge & \beta = \cap \\ \neg & \beta = \overline{} \\ 0^m & \beta = \emptyset \\ 1^m & \beta = \mathcal{B}^*, \end{cases} \quad (17)$$

where \vee , \wedge , and \neg are respectively bitwise or, and, and negation.

This homomorphism is *one-way* for two independent reasons. First, F is not injective, i.e., multiple values in \mathcal{B}^* (countably infinite many, since \mathcal{B}^* is countably infinite) must map to at least one value in \mathcal{B}^m . Second, the hash function h is a cryptographic hash function that is resistant to preimage attacks, i.e., given a y it is *hard* to find a x such that $h(x) = y$ but given an x it is *easy* to find a y such that $h(x) = y$.

Thus, given a set $\{x, y, z\} = \{x\} \cup \{y\} \cup \{z\}$, F maps this to $h(x) \vee h(y) \vee h(z)$.

These model *positive approximate sets* with a false positive rate given by the following theorem.

Theorem 3.1. The false positive rate is a function of m and n and is given by

$$\varepsilon(m, n) = \alpha^m(n) \quad (18)$$

where

$$\alpha(n) := \left(1 - 2^{-(n+1)} \right). \quad (19)$$

Proof. Proof here. See paper sheets that are star-numbered. \square

The *bit-rate*, the bits per element, is given by the following corollary.

Corollary 3.1.1.

$$b(n, \varepsilon) = \frac{\log_2 \varepsilon}{n\alpha(n)} \quad (20)$$

Proof. Proof here. See paper sheets that are star-numbered. \square

The *partial derivative* of b with respect to ε is given by

$$\frac{\partial b}{\partial \varepsilon} = -\frac{b}{B(\varepsilon)}, \quad (21)$$

where $B(\varepsilon) := -\varepsilon \log_2 \varepsilon$ is the *binary entropy function*. This partial is less than zero everywhere but finds its minimum at $\frac{1}{e}$.

The *absolute efficiency* of this homomorphism, when implemented as a data structure for a false positive rate ε is given by

$$E(n) = -n \log_2 \alpha(n), \quad (22)$$

which is exponential with respect to n .

As $n \rightarrow \infty$, $E(n)$ goes to 0. It is only practical for small values of n , where n is the number of elements of the set.

The *identity* relation has false positives given by the following.

Theorem 3.2. *False positives on the equality predicate occurs with probability*

$$\varepsilon(m) = ?, \quad (23)$$

where m is the bit length of the hash function in the Boolean cipher algebra.

Proof. Two sets are said to be identical if they map to the same bit string, and thus, as before, *false positives* are possible,

$$P[GX = GY \mid X \neq Y] \quad (a)$$

which occurs with probability ?... \square

TODO: figured out membership fpr, now figure out subset fpr, identity fpr also.

4 Cipher Boolean ring

Say we have the Boolean ring

$$A := (2^{\mathcal{B}^*}, \Delta, \cap, \overline{}, \emptyset) \quad (24)$$

where \mathcal{B}^* is the free semigroup over \mathcal{B} , which is the closure over concatenation $\vdash: \mathcal{B}^* \times \mathcal{B}^* \mapsto \mathcal{B}^*$.

Consider the Boolean ring

$$B := (\mathcal{B}^m, \oplus, \wedge, \text{id}, 0^m) \quad (25)$$

and suppose we wish to provide a homomorphism of type $A \mapsto B$ that *only* preserves identities.

Definition 4.1. The homomorphism $G: A \mapsto C$ is defined as

$$G(\beta) := \begin{cases} h(x) & \beta \in \mathcal{B}^* \\ \oplus & \beta = \Delta \\ \wedge & \beta = \cap \\ \text{id} & \beta = \overline{} \\ 0^m & \beta = \emptyset \end{cases} \quad (26)$$

where \oplus , \wedge , and id are respectively bitwise exclusive-or, bitwise and, and identity.

This homomorphism is *one-way* for the same reason as before. The identity relation

$$(G \mathcal{A})(F \Delta)(G \mathcal{B}) = G(\mathcal{A} \Delta \mathcal{B}) \quad (27)$$

is guaranteed by the homomorphism.²

There is an *isomorphism* between G and F since $\mathcal{A} \cup \mathcal{B} := (\mathcal{A} \Delta \mathcal{B}) \Delta (\mathcal{A} \cap \mathcal{B})$. If the trusted system applies only the *group*

$$D := (\mathcal{B}^m, \oplus, \text{id}, 0^m) \quad (28)$$

the number of 1's and the cardinality of the set has no relation.

Given a set $\{x, y, z\} = \{x\} \cup \{y\} \cup \{z\}$, G maps this to $h(x) \oplus h(y) \oplus h(z)$.

Since the hash function h is a random oracle, by the property of the exclusive-or operation each set in A maps to a *random* bit string in \mathcal{B}^* . That is to say, it is indistinguishable from random noise. However, it only supports *identity* relations.

Theorem 4.1. False positives on the equality predicate occurs with probability

$$\varepsilon(m) = 2^{-m}, \quad (29)$$

where m is the bit length of the hash function in the Boolean cipher ring.

Proof. Two sets are said to be identical if they map to the same bit string, and thus, as before, *false positives* are possible,

$$P[G \mathcal{X} = G \mathcal{Y} \mid \mathcal{X} \neq \mathcal{Y}] \quad (a)$$

which occurs with probability 2^{-m} since each set maps to a random bit string of size m . \square

5 The Singular Hash Set

We would like to have properties of both homomorphisms F and G but with a much smaller space complexity.

In what follows, we provide a theoretical implementation of the *cipher set* that obtains optimality in the following ways:

1. The space complexity obtains the theoretical lower-bound of a random approximate positive set with an expected false positive rate ε .
2. The entropy obtains the upper bound for the given expected space complexity.

²Observe that if \mathcal{A} and \mathcal{B} are disjoint, $\mathcal{A} \Delta \mathcal{B} = \mathcal{A} \cup \mathcal{B}$.

Definition 5.1. The data type for the cryptographic singular hash set is defined as $\text{SHS}_{\mathcal{F}\mathcal{X}} := \mathcal{B}^k \times \mathcal{B}^*$ with a value constructor $\text{shs}_{\mathcal{F}\mathcal{X}}: 2^{\mathcal{X}} \times \{2^{-k} \in \mathbb{Q} \mid k \in \mathbb{N}\} \mapsto \text{SHS}_{\mathcal{F}\mathcal{X}}$ defined as

$$\text{shs}_{\mathcal{F}\mathcal{X}}(\mathcal{A}, \varepsilon) := \langle q, n' \rangle \quad (30)$$

where

$$\begin{aligned} \mathcal{A} &= \{x_{(1)}, \dots, x_{(m)}\}, \\ k &:= -\log_2 \varepsilon, \\ h_k(b, n) &:= \text{tr} \left(h(b \# n'), k \right), \\ \phi(w, n) &:= h_k(E(w), n), \\ \Omega_l &:= \{ \phi(x, l) \in \mathcal{B}^k \mid w \in F(\mathcal{A}) \}, \\ n &:= \min \{ j \in \mathbb{N} \mid \Omega_j \in 2^{\mathcal{B}^k} \wedge |\Omega_j| = 1 \}, \\ q &:= \phi(F x_{(1)}, n). \end{aligned} \quad (31)$$

where $F: \mathcal{X} \mapsto \mathcal{W}$ is a homomorphism and $E: \mathcal{W} \mapsto \mathcal{B}^*$ is a prefix-free coder.³

If F is the *identify* function $\text{id}: \mathcal{X} \mapsto \mathcal{X}$ then the singular hash set models a random approximate positive set over \mathcal{X} , denoted by $(2^{\mathcal{X}})^+$. However, typically, the singular hash set is intended to be used for as a *cipher set*, a parametric type $\boxed{\check{\mathcal{X}}}^+$ where $F: \mathcal{X} \mapsto \check{\mathcal{X}}$, e.g., $x \mapsto \check{x}$ where \check{x} may be, say, an *equivalence class* $\{\check{x}_1, \dots, \check{x}_k\}$ in which each of the elements map back to x under $F^{-1}: \check{\mathcal{X}} \mapsto \mathcal{X}$.

Theorem 5.1. The data type SHS with value constructor $\text{shs}: \mathcal{P}(\mathcal{B}^*) \times \{2^{-k} \in \mathbb{Q} \mid k \in \mathbb{N}\} \mapsto \text{SHS}$ over the computational basis $\in: \text{SHS} \times \mathcal{B}^* \mapsto \mathcal{B}$ is a random approximate positive set of $\check{\mathcal{A}}$ with a false positive rate ε over $\check{\mathcal{X}} \setminus \check{\mathcal{A}}$. That is, a priori,

$$\text{shs}_{\mathcal{X}}(\mathcal{A}, \varepsilon) \sim \check{\mathcal{A}}_{\varepsilon}^+ \quad (32)$$

for any $\mathcal{A} \in 2^{\mathcal{X}}$ and $\varepsilon \in \{2^{-k} \in \mathbb{Q} \mid k \in \mathbb{N}\}$.

Proof. In order for the value $\text{shs}_{\mathcal{X}}(\mathcal{A}, \varepsilon)$ to be a random approximate set with a distribution given by $\check{\mathcal{A}}_{\varepsilon}^+$, it must satisfy two conditions as specified by ??:

1. $\check{\mathcal{A}}$ is a subset of $\check{\mathcal{A}}^+$. This condition guarantees that no *false negatives* may occur.
2. An element in $\check{\mathcal{X}}$ that is not a member of $\check{\mathcal{A}}$ is a member of $\check{\mathcal{A}}^+$ with a probability ε , denoted the false positive rate, i.e.,

$$\text{P}[x \in \check{\mathcal{A}}^{\pm} \mid x \notin \check{\mathcal{A}}] = \varepsilon \quad (\text{a})$$

for any $x \in \mathcal{X}$.

To prove the first condition, note that ?? tests any element x for membership in \mathcal{S}^+ by computing the hash of x concatenated with the bit string b_n and returning **true** if the hash is h_k where definition 5.1 finds bit strings b_n and h_k such that each element of \mathcal{S} concatenated with b_n hashes to h_k .

To prove the second condition, suppose we have a set $\mathcal{S} = \{x_1, \dots, x_m\}$ and each element in \mathcal{S} hashes to $y = h(x_1)$. By ??, $h: \mathcal{B}^* \mapsto \mathcal{B}^k$ approximates a random oracle and thus uniformly distributes over its domain of 2^k possibilities. Since y is a particular element in \mathcal{B}^k , the probability that an element not in \mathcal{S} hashes to y is 2^{-k} . \square

³If E is not prefix-free, then an approximation error independent of parameter ε is introduced.

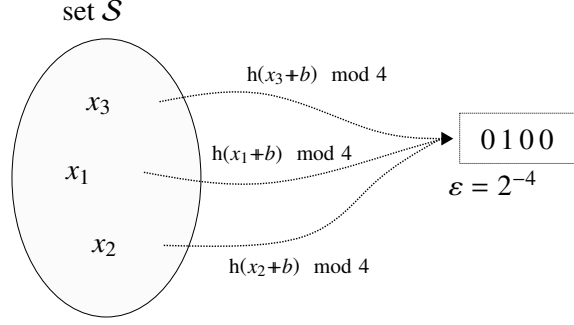


Figure 1: *Singular Hash Set* over a *countably infinite* universe

In particular, we consider a data type that supports both the membership and identity relations while obtaining the information-theoretic expected lower-bound for space with respect to the false positive rate. Since it obtains the information theoretic lower-bound, the bit strings generated naturally obtain *maximum entropy*, i.e., any set $\mathcal{A} \in 2^{\mathcal{X}}$ a priori maps to random bit string in its representation. The only information it leaks is related to its expected bits per element, i.e., it is possible to estimate the cardinality of the set with high certainty.⁴

The singular hash set obtains the theoretically optimal lower bound for a random approximate set (and by definition therefore obtains maximum entropy) while providing for set-membership.

The singular hash set is a data type that implements the *oblivious* random positive approximate set abstract data type. The implementation consists of a product data structure (tuple) $\mathcal{B}^k \times \mathcal{B}^*$, an algorithm that *generates* the data structure, and an algorithm that implements the *member-of* function by appropriately *querying* the data structure.

The *membership* relation is defined by the predicate $\in: \check{\mathcal{X}} \times \text{SHS}_{\mathcal{X}} \mapsto \mathcal{B}$, is defined as

$$\langle h, b \rangle \in \check{x} := q' \bmod N \quad (33)$$

where

$$\begin{aligned} k &:= \lceil \log_2 N \rceil, \\ r &:= h(x' + b), \\ q &:= \text{tr}(r, k). \end{aligned} \quad (34)$$

Definition 5.2.

$$\in: \mathcal{X} \times \text{SHS}_{\mathcal{X}} \mapsto \mathcal{B} \quad (35)$$

5.1 Time and space complexity

The probability that every element of \mathcal{S} collides for a particular bit string in ?? is given by the following theorem.

Theorem 5.2. *The number of time steps the algorithm for $\text{shs}(m, \epsilon)$ is geometrically distributed*

$$Q \sim \text{GEO}(\epsilon^{m-1}) \quad (36)$$

with an expected number of steps given by

$$\epsilon^{-(m-1)}. \quad (37)$$

⁴Of course, the objective set may be poisoned with nonsense values before generating the set such that, for instance, every set is expected to have the same bit length.

Proof. Suppose we have a set of bit strings $\mathcal{A} = \{b_1, \dots, b_m\}$ and b_1 hashes to $y = h_k(b_1, n)$ where $h_k: \mathcal{B}^* \mapsto \mathcal{B}^k$ is a random hash function that uniformly distributes over its domain of 2^k possibilities. Since y is a particular element in \mathcal{B}^k , the probability that b_j for $j = 2, \dots, m$ hashes to y is given by

$$\frac{1}{2^k} = \varepsilon. \quad (a)$$

Since h_k is a random hash function, the hashes of b_1, \dots, b_m are independent. Thus, the joint probability that b_2, \dots, b_m hash to y is given by the product of their marginal probabilities

$$\varepsilon^{m-1}. \quad (b)$$

□

Theorem 5.3. *The expected bit length of the Singular Hash Set obtains the information-theoretic lower-bound given by*

$$-\log_2 \varepsilon \text{ bits/element}, \quad (38)$$

where ε is the false positive rate.

Proof. Suppose $|\mathcal{A}| = m$. By definition 5.1, $\text{shs}_{\text{FX}}(\mathcal{A}, \varepsilon)$ is a value constructor that returns a pair of elements, $\langle n', q \rangle$ which has a bit length $\ell(n') + \ell(q)$. The bit length of $\ell(q)$ is just $-\log_2 \varepsilon$. The bit length of n' depends on the minimum value of n found.

Computationally, we search in the order of increasing n , from 0 to ∞ . The first case when a perfect collision occurs is a geometric distribution $Q \sim \text{GEO}(p = \varepsilon^{m-1})$. By definition 2.2, the n -th trial uniquely maps to a bit string of length $\lfloor \log_2 n \rfloor$. Thus, the bit string has a random length given approximately by $N = \log_2 Q$ bits. We approximate the logarithm with a second-order Taylor series around the *expected* value of Q as given by

$$N \approx \log_2 E[Q] - \frac{(Q - E[Q])^2}{2 E[Q]^2} \log_2 e \text{ bits}. \quad (a)$$

We are interested in the *expected* value of N ,

$$E[N] \approx \log_2 E[Q] - \frac{V[Q]}{2 E[Q]^2} \log_2 e \text{ bits}. \quad (b)$$

The expectation and variance of Q is known to be $1/p$ and $(1-p)/p^2$ respectively, and thus we may rewrite eq. (b) as

$$E[N] \approx -\log_2 p - \frac{1-p}{2} \log_2 e \text{ bits}. \quad (c)$$

Substituting ε^{m-1} for p in eq. (c) results in

$$E[N] \approx -(m-1) \log_2 \varepsilon - \frac{1 - \varepsilon^{m-1}}{2} \log_2 e \text{ bits}. \quad (d)$$

In total, the bit length of $\langle n', q \rangle$ is thus

$$\ell = -m \log_2 \varepsilon + \frac{1 - \varepsilon^{m-1}}{2} \log_2 e \text{ bits} \quad (e)$$

which asymptotically converges to $-\log_2 \varepsilon$ bits per element as $m \rightarrow \infty$. □

By ??, ?? has an expected time complexity that grows exponentially as m grows. The algorithm is intended to illustrate theoretical properties, not necessarily be used in practice.

5.2 Probability distributions

The bit representation of a singular hash set of \mathcal{A} is uncorrelated with the specific elements of but is highly correlated with the its *cardinality* $|\mathcal{A}|$.

The probability that $B = b$ is given by

$$f_{B|N}(b | n) = 2^{-n} \mathbb{1}_{\mathcal{B}^n}(b). \quad (39)$$

Proof. Sketch proof. □

Given a false positive rate ε , the probability that $H = h$ is given by

$$f_H(h | \varepsilon) = \varepsilon \mathbb{1}_{\mathcal{B}^k}(h), \quad (40)$$

where $k = -\log_2 \varepsilon$.

Theorem 5.4. *The random bit length N has a probability mass function given by*

$$f_N(n | m, \varepsilon) = \left(\frac{1}{q} - q \right) q^{2^n} \mathbb{1}_N(n), \quad (41)$$

where $q = 1 - \varepsilon^{m-1}$, m is the cardinality of the objective set, and ε is the false positive rate.

Proof. The probability that a bit string of length n occurs is given by the probability that Q realizes some value in the range $\{2^n, \dots, 2^{n+1} - 1\}$, which is given by

$$P[2^n - 1 < Q \leq 2^{n+1} - 1] = F_Q(2^{n+1} - 1) - F_Q(2^n) \quad (a)$$

$$= \left(1 - q^{2^{n+1}-1} \right) - \left(1 - q^{2^n-1} \right) \quad (b)$$

$$= q^{2^n-1} - q^{2^{n+1}-1}, \quad (c)$$

which is equivalent to the result. □

Given a false positive rate $\mathcal{E} = \varepsilon$ and an objective set of cardinality m , the joint distribution of B , H , and N is given by

$$f_{B,H,N}(b, h, n | m, \varepsilon) = 2^{-n} \varepsilon f_N(n | m, \varepsilon) \mathbb{1}_{\mathcal{B}^n \times \mathcal{B}^k \times \mathbb{N}}(b, h, n) \quad (42)$$

where $k = -\log_2 \varepsilon$.

In the singular hash set, the false positive rate and bit length N are *observable*. Therefore, a primary statistic of interest is the joint distribution of B and H given $N = n$ and $\mathcal{E} = \varepsilon$, which is given by

$$f_{B,H|N}(b, h | n, \varepsilon) = 2^{-n} \varepsilon \mathbb{1}_{\mathcal{B}^n \times \mathcal{B}^k}(b, h) \quad (43)$$

where $k = -\log_2 \varepsilon$.

The marginal distribution of B is given by marginalizing over the joint distribution with respect to B , which yields the result

$$f_B(b | \varepsilon, m) = 2^{-\ell(b)} f_N(\ell(b) | m, \varepsilon). \quad (44)$$

A few key features of B is that it realizes the empty string with probability ε^{m-1} and as n goes to infinity, the probability that $\ell(B) = n$ goes to 0.

NOTE: the singular hash set does not expose correlations by bit-wise operations, which may be a good thing, unlike in the bitwise models discussed previously.

The entropy of N is given by the following theorem.

Theorem 5.5. ?

Proof.

$$\mathcal{H}(N) = -E[\log_2 f_N(N)]. \quad (a)$$

$$\mathcal{H}(N) = -\sum_{n=0}^{\infty} \log_2 f_N(n) f_N(n) \quad (b)$$

$$= -\sum_{n=0}^{\infty} \log_2 \left(q^{2^n-1} (1 - q^{2^n}) \right) f_N(n). \quad (c)$$

$$\begin{aligned} \mathcal{H}(N) = & -\log_2 q \sum_{n=0}^{\infty} (2^n - 1) f_N(n) - \\ & \sum_{n=0}^{\infty} \log_2 (1 - q^{2^n}) f_N(n). \end{aligned} \quad (d)$$

$$\begin{aligned} \mathcal{H}(N) = & -\log_2 q \left[\sum_{n=0}^{\infty} (2^n f_N(n)) \right] - \\ & \sum_{n=0}^{\infty} \log_2 (1 - q^{2^n}) f_N(n). \end{aligned} \quad (e)$$

□

Theorem 5.6. *The random bit string that codes the singular hash set is a maximum entropy coder for the approximate set abstract data type.*

Proof. The result immediately follows from the fact that the bit string is *incompressible* and thus obtains maximum entropy. The bit string $h_k \in \mathcal{B}^k$ is, a priori, a random bit string uniformly distributed over \mathcal{B}^k by the property of the cryptographic hash function. The bit string $b_n \in \mathcal{B}^*$ realizes a length n with probability $f_N(n)$.⁵ Given $N = n$, the bit string $b_n \in \mathcal{B}^n$ is, a priori, a random bit string uniformly distributed over \mathcal{B}^n . □

There are two predictable regularities in the bit string representation of a singular hash set.

1. The *random bit length* N has relatively low entropy and is expected to obtain the information-theoretic lower-bound.
2. The false positive rate has no uncertainty, i.e., zero entropy.

By construction, we are able to estimate that an element is a member or non-member of \mathcal{A} by determining if it is a member or non-member of $\hat{\mathcal{A}}^+$, e.g., positive predictive value and other binary classification measures.

By construction, the singular hash set obtain the information-theoretic lower-bound, which may be used to estimate the cardinality of the objective sets being modeled by a singular hash set. For instance, given a set \mathcal{A} , the random bit length N of bit string b_n has a probability mass concentrated around the theoretical lower-bound. Therefore, a *method-of-moments* estimator of the cardinality of the \mathcal{A} is given by

$$|\hat{\mathcal{A}}| = -\frac{\ell(\mathcal{A}^+)}{\log_2 \varepsilon}, \quad (45)$$

where ℓ is the bit length function and ε is the *false positive rate*.

⁵As such, we can further compress b_n using a geometric coder that assigns shorter bit strings to more probable lengths, but we require an *in-place* bit string and so stick with the original code.

5.3 Relaxing optimality

By ??, the time complexity is *exponential* with respect to the cardinality m of the objective set being modeled with a cipher set. We can trade space complexity for time complexity to design more practical algorithms.

Let \mathcal{A} be the objective set. A general approach is to use a *multi-level* hashing scheme such that each level generates smaller independent sub-problems.

We consider a two-level scheme. At the first level, we find a hash function $h_k: \mathcal{X} \mapsto \{0, 1, \dots, k-1\}$ that maps m/k elements in \mathcal{A} to each element in the codomain.

Definition 5.3. *The cryptographic α -perfect hash function is a variation of the perfect hash function. For each element in the codomain, proportion α of the domain maps to it.*

There are two degenerate cases to consider. Suppose \mathcal{A} is the objective set and $|\mathcal{A}| = m$. If $\alpha = 1$, then it reduces to a *singular* hash function where every element of \mathcal{A} maps to the same value. If $\alpha = \frac{1}{m}$, then it reduces to a *minimal perfect hash function* of \mathcal{A} where the hash function restricted to \mathcal{A} is a bijection between \mathcal{A} and $\{0, \dots, m-1\}$.

The cryptographic m/k -perfect hash function has a space and time complexity given by the following theorem.

Theorem 5.7. *A cryptographic m/k -perfect hash function has an expected space complexity given by*

$$- \alpha \log_2 \alpha \quad (46)$$

and an expected time complexity given by

$$\sqrt{\frac{k}{\alpha}} \alpha^{\frac{k}{2}} \quad (47)$$

where $\alpha := \frac{m}{k}$.

Proof. Consider the family of hash functions of the type $\mathcal{X} \mapsto \{1, 2, \dots, k\}$. When we restrict this family to some objective set \mathcal{A} of cardinality m , we have hash functions of type $\mathcal{A} \mapsto \{1, 2, \dots, k\}$. There are exactly k^m hash functions of this type.

We are interested only in those hash functions of this type that evenly k -partition the domain into the codomain, i.e., the preimage of each element in the codomain consists of $\alpha := m/k$ elements in the domain.

To count the number of possible functions, we choose α of the m elements in the domain for element 1 in the codomain. There are a total of $\binom{m}{\alpha}$ ways of making this selection. We remove these α elements from the domain. We repeat the procedure again. We choose α of the remaining $m - \alpha$ elements in the domain for element 2 in the codomain. There are a total of $\binom{m-\alpha}{\alpha}$ ways of making this selection. If there are a ways to make a first choice and b ways to make a second choice, by the product rule there are ab ways to make those two choices. Thus, continuing this pattern, there are

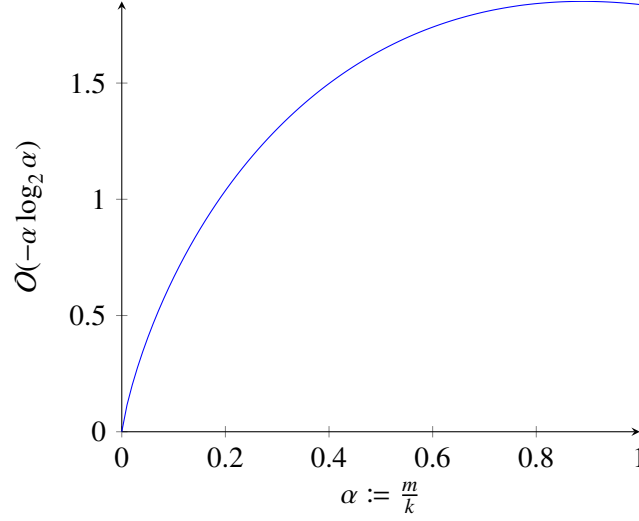
$$\prod_{j=1}^k \binom{m - j\alpha}{\alpha} \quad (a)$$

ways to make the selections for the k elements in the codomain, which simplifies to

$$\frac{m!}{(\alpha!)^k}. \quad (b)$$

To generate perfect hash functions, we append some bit string to each input. By the property of cryptographic hash functions, each bit string serves as an index into the cryptographic hash function family

Figure 2: The expected bits per element of the cryptographic k -perfect hash function.



$\mathcal{A} \mapsto 0, 1, \dots, k-1$. Thus, since there are k^m functions in this family in total, and only $\frac{m!}{(\alpha!)^k}$ m/k -perfect hash functions of this type, each time we index into the family, with probability

$$p \sim e \sqrt{\frac{k}{\alpha}} \alpha^{\frac{k}{2}} \quad (\text{c})$$

will the hash function be a α -perfect hash function.

The probability of successfully finding an m/k -perfect hash function for each trial bit string is thus *geometrically* distributed with a probability of success p .

The time complexity is just the *expected* number of trials, which is $1/p$ for the geometric distribution, and the length of the successful bit string is just $\log_2 1/p$. \square

We employ the m/k -perfect hash function to construct a more practical variation of the singular hash set.

We k -partition the m elements in a specified set using the m/k -perfect hash function, giving us m/k elements per bin, and then we apply the singular hash set to each of these independently.

This is a negative binomial distribution.

Suppose we have two cryptographic hash functions, $h_1: \mathcal{B}^* \mapsto \mathcal{B}^w$ and $h_2: \mathcal{B}^w \mapsto \mathcal{B}^k$ where, in general, $w < k$ but this is not necessarily the case. If $w \gg |\mathcal{A}|$, then the size of the cipher set is just $-w \log_2 \varepsilon$ where $\varepsilon = 2^{-k}$ since each element of \mathcal{A} , with high probability, maps to a unique hash by h_1 .

The first-level hash function $h_k := \mathcal{X} \mapsto \{0, 1, \dots, k-1\}$ is defined as

$$h_k(x) := \langle h, n' \rangle \quad (48)$$

where

$$\begin{aligned} \phi(w, n) &:= h(n' + x') \mod k, \\ \mathcal{Y}_l &:= \{ \phi(x, l) \in \mathcal{B}^k \mid x \in \mathcal{A} \}, \\ n &:= \min \left\{ \{ j \in \mathbb{N} \mid \mathcal{Y}_j \in 2^{\mathcal{B}^k} \wedge |\mathcal{Y}_j| = 1 \} \right\}. \end{aligned} \quad (49)$$

The first application of the cryptographic hash function maps each element to a particular index. By the property of the cryptographic hash function, the elements are expected to be equally partitioned into

the k indices. Next, at each index, we associate an independent bit string that serves the same purpose as the bit string b_n in the singular hash set. We then independently apply the singular hash set algorithm to each partition. Let the number of trials necessary for singular hash set of the j -th partition be denoted by T_j . The total expected number of trials is just the sum of the geometrically distributed random variables $T = T_1 + T_2 + \dots + T_k$, which is a *negative binomial distribution*,

$$T \sim \text{NB} \left(k, \varepsilon^{\frac{m}{k}-1} \right). \quad (50)$$

The extra degree of freedom we have with the number of partitions k means we can quantitatively trade space complexity for time complexity.

Theorem 5.8. *The value of k that is expected to yield a solution at t trials is given by*

$$k = \alpha E[N], \quad (51)$$

where $\frac{1}{\alpha} = \log_2 \left(\frac{\varepsilon t}{k+t} \right)$.

Proof. We take the approach of solving for some t number of trials and finding a lower-bound k that satisfies the equation,

$$E[T] = t. \quad (a)$$

Substituting the expectation of the negative binomial into the above yields the equation

$$\frac{k \varepsilon^{m/k-1}}{1 - \varepsilon^{m/k-1}} = t. \quad (b)$$

Solving for k yields the result. \square

We know that, approximately, the bit length is given by

$$N = \log_2 T, \quad (52)$$

which may be solved as described in ??.

We may also quantify the probability that, after t trials, some jointly chosen value of k , t , and ε fails to realize a solution. In fact, we can quantify the probability that r of the m elements fails, which yields a *rate-distortion* $\frac{r}{m}$.

If after t trials no solution is found, we may either continue with the search or stop short, resulting in not just an approximate set with *false positives* but also *false negatives*. We refer to this outcome as a *rate-distorted set*, where the rate distortion is a function of time. However, note that we could also make the rate-distortion a function of *space*, i.e., limit the maximum size of the bit string code and choose the code that yields the smallest false negative rate.

If we are interested in rate-distorted sets with false negative rates greater than 0, then an alternative approach is to consider the probability of failure for a particular k , ε , and m . By the survival function of the negative binomial, the probability that more than t trials (or $\log_2 t$ bits) is required for a particular parameterization k and $p = \varepsilon^{\frac{m}{k}-1}$ is given by

$$I_p(t+1, k), \quad (53)$$

where I_p is the *regularized incomplete beta function*.

Theorem 5.9. *A value of type SHS_X^k constructed with $\text{shs}_X(\mathcal{A}, \varepsilon \mid k = m)$ models a cryptographic perfect hash function $h_{\mathcal{A}}: X \mapsto \{0, 1, \dots, -\log_2 \varepsilon\}$ where $\mathcal{A} \subseteq X$ and $r = \frac{|\mathcal{A}|}{-\log_2 \varepsilon}$.*

5.4 Boolean algebra

The data structure does not trivially model a Boolean algebra with simple bit-wise operators on its representation, as is the case with F and G, but it may be *lifted* to a Boolean algebra by the fact that

$$\mathcal{A} \cup \mathcal{B} = \{x \in \mathcal{X} \mid x \in \mathcal{A} \vee x \in \mathcal{B}\} \quad (54)$$

and other similar identities. Specifically, we define a *recursive* type for the Boolean algebra as

$$\mathbf{B}_{\tilde{\mathcal{X}}} := 0_{\tilde{\mathcal{X}}} + 1_{\tilde{\mathcal{X}}} \quad (55)$$

$$+ \text{Just SHS}_{\tilde{\mathcal{X}}} \quad (56)$$

$$+ \text{Join } \mathbf{B}_{\tilde{\mathcal{X}}} \mathbf{B}_{\tilde{\mathcal{X}}} \quad (57)$$

$$+ \text{Meet } \mathbf{B}_{\tilde{\mathcal{X}}} \mathbf{B}_{\tilde{\mathcal{X}}} \quad (58)$$

$$+ \text{Not } \mathbf{B}_{\tilde{\mathcal{X}}} . \quad (59)$$

Then, for instance, set-union $\cup: \mathbf{B}_{\tilde{\mathcal{X}}} \times \mathbf{B}_{\tilde{\mathcal{X}}} \mapsto \mathbf{B}_{\tilde{\mathcal{X}}}$ is defined as $a \cup b := \text{Join } a b$.

Over the Boolean algebra $\mathbf{B}_{\tilde{\mathcal{X}}}$, the membership predicate

$$\in: \mathbf{B}_{\tilde{\mathcal{X}}} \times \mathbf{B}_{\tilde{\mathcal{X}}} \mapsto \mathcal{B} \quad (60)$$

is defined as

$$x \in 1_{\tilde{\mathcal{X}}} := 1 , \quad (61)$$

$$x \in 0_{\tilde{\mathcal{X}}} := 0 , \quad (62)$$

$$x \in \text{Just } a := x \in a , \quad (63)$$

$$x \in \text{Join } a b := (x \in a) \vee (x \in b) , \quad (64)$$

$$x \in \text{Meet } a b := (x \in a) \wedge (x \in b) , \quad (65)$$

$$x \in \text{Not } a := \neg(x \in a) . \quad (66)$$

The nature of the composition reveals something about the sets and decreases entropy.

With respect to the above, any Boolean algebra

$$A := \left(2^{\mathcal{X}}, \cup, \cap, \overline{}, \emptyset, \mathcal{X} \right) . \quad (67)$$

may be mapped to the Boolean algebra

$$B := \left(\mathbf{B}_{\tilde{\mathcal{X}}}, \cup, \cap, \overline{}, 0_{\tilde{\mathcal{X}}}, 1_{\tilde{\mathcal{X}}} \right) . \quad (68)$$

Interestingly, the mapping from A to B is not a homomorphism since, for instance, $\text{shs}_{\tilde{\mathcal{X}}}(\mathcal{A}) \cup \text{shs}_{\tilde{\mathcal{X}}}(\mathcal{B})$ does not necessarily equal $\text{shs}_{\tilde{\mathcal{X}}}(\mathcal{A} \cup \mathcal{B})$ due to the random approximate set model. That is, the mapping only *probabilistically* preserves relations according to the *higher-order* random approximate set model. One might say that it is *asymptotically* homomorphic ($\varepsilon \rightarrow 0$).

Rate-distorted cipher set The *false positive rate* of an approximate cipher set is a *rate-distortion*. The rate distortion is a function of strictly space constraints.

We consider another kind of rate distortion that generates *false negatives* which is either a function of space or time constraints.

The result is an approximate cipher set with both a first-order false positive *and* false negative rates. Note that cipher sets with higher-order false negatives and false positives may also be generated by *composing* the

cipher sets, but here we consider a value constructor of cipher sets that generates first-order approximation errors.

The number of elements that collide on the *singular hash* on any particular trial is binomially distributed,

$$T \sim B(m - 1, p) \tag{69}$$

where $p = \varepsilon$. When we are only interested in m ‘successes’, the probability is given as before, ε^{m-1} .

There are a few different ways of treating the rate distortion ω and ε .

We consider the false positive rate ε first. Generally, we specify a value that is a compromise between space and accuracy. Since the cost of a false positive and the cost of a false negative may differ, different values for each may be preferable.

One way is to treat α as a random variable that is induced by a maximum bit rate per element, e.g., on average the maximum number of bits per element is b . Thus, if there are m elements, the constructor is limited to finding a bit string of length bm .

Another way to proceed is to say that we are satisfied with proportion $(1 - \alpha)$ elements or more being perfectly hashed.