

A preliminary version of this paper appears in *Advances in Cryptology - EUROCRYPT 2009, 28th Annual International Cryptology Conference*, A. Joux ed., LNCS, Springer, 2009.

## Order-Preserving Symmetric Encryption

ALEXANDRA BOLDYREVA\*    NATHAN CHENETTE<sup>†</sup>    YOUNHO LEE<sup>‡</sup>  
ADAM O'NEILL<sup>§</sup>

November 4, 2012

### Abstract

We initiate the cryptographic study of order-preserving symmetric encryption (OPE), a primitive suggested in the database community by Agrawal et al. (SIGMOD '04) for allowing efficient range queries on encrypted data. Interestingly, we first show that a straightforward relaxation of standard security notions for encryption such as indistinguishability against chosen-plaintext attack (IND-CPA) is unachievable by a practical OPE scheme. Instead, we propose a security notion in the spirit of pseudorandom functions (PRFs) and related primitives asking that an OPE scheme look “as-random-as-possible” subject to the order-preserving constraint. We then design an efficient OPE scheme and prove its security under our notion based on pseudorandomness of an underlying blockcipher. Our construction is based on a natural relation we uncover between a random order-preserving function and the hypergeometric probability distribution. In particular, it makes black-box use of an efficient sampling algorithm for the latter.

---

\*School of Computer Science, Georgia Institute of Technology, 266 Ferst Drive, Atlanta, GA 30332, USA. E-mail: [sasha@gatech.edu](mailto:sasha@gatech.edu).

<sup>†</sup>Department of Mathematical Sciences, Clemson University, O-110 Martin Hall, Box 340975, Clemson, SC 29634, USA. E-mail: [nchenet@clemson.edu](mailto:nchenet@clemson.edu). Most of the work done while at the Georgia Institute of Technology.

<sup>‡</sup>Department of Information and Communication Engineering Yeungnam University, Republic of Korea. E-mail: [younholee@yu.ac.kr](mailto:younholee@yu.ac.kr). Work done while at the Georgia Institute of Technology.

<sup>§</sup>Department of Computer Science, Boston University, 111 Cummington St., Boston, MA 02215, USA. E-mail: [amoneill@bu.edu](mailto:amoneill@bu.edu). Work done while at the Georgia Institute of Technology and University of Texas.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Preliminaries</b>	<b>6</b>
<b>3</b>	<b>Order-Preserving Encryption and its Security</b>	<b>8</b>
3.1	Order-Preserving Encryption (OPE) . . . . .	8
3.2	Security of OPE . . . . .	8
<b>4</b>	<b>Lazy Sampling a Random Order-Preserving Function</b>	<b>11</b>
4.1	The Hypergeometric Connection . . . . .	11
4.2	The <b>LazySample</b> Algorithms . . . . .	12
4.3	Correctness . . . . .	13
4.4	Efficiency . . . . .	15
4.5	Realizing HG . . . . .	16
<b>5</b>	<b>Our OPE Scheme and its Analysis</b>	<b>17</b>
5.1	The TapeGen PRF . . . . .	17
5.2	Our OPE Scheme and its Analysis . . . . .	20
5.3	On Choosing $N$ . . . . .	21
<b>6</b>	<b>On Using the Negative Hypergeometric Distribution</b>	<b>22</b>
6.1	Construction of the NHG-based OPE Scheme . . . . .	22
6.2	Correctness . . . . .	22
6.3	Efficiency of the NHG Scheme . . . . .	25

# 1 Introduction

**MOTIVATION.** Order-preserving symmetric encryption (OPE) is a deterministic encryption scheme (aka. cipher) whose encryption function preserves numerical ordering of the plaintexts. OPE has a long history in the form of one-part codes, which are lists of plaintexts and the corresponding ciphertexts, both arranged in alphabetical or numerical order so only a single copy is required for efficient encryption and decryption. One-part codes were used, for example, during World War I [3]. A more formal treatment of the concept of order-preserving symmetric encryption (OPE) was proposed in the database community by Agrawal et al. [1]. The reason for new interest in such schemes is that they allow efficient range queries on encrypted data. That is, a remote untrusted database server is able to index the (sensitive) data it receives, in encrypted form, in a data structure that permits efficient range queries (asking the server to return ciphertexts in the database whose decryptions fall within a given range, say  $[a, b]$ ). By “efficient” we mean in time logarithmic (or at least sub-linear) in the size of the database, as performing linear work on each query is prohibitively slow in practice for large databases.

In fact, OPE not only allows efficient range queries, but allows indexing and query processing to be done exactly and as efficiently as for unencrypted data, since a query just consists of the encryptions of  $a$  and  $b$  and the server can locate the desired ciphertexts in logarithmic-time via standard tree-based data structures. Indeed, subsequent to its publication, [1] has been referenced widely in the database community, and OPE has also been suggested for use in in-network aggregation on encrypted data in sensor networks [32] and as a tool for applying signal processing techniques to multimedia content protection [14]. Yet a cryptographic study of OPE in the provable-security tradition never appeared. Our work aims to begin to remedy this situation.

**RELATED WORK.** Our work extends a recent line of research in the cryptographic community addressing efficient (sub-linear time) search on encrypted data, which has been addressed by [2] in the symmetric-key setting and [6, 11, 7] in the public-key setting. However, these works focus mainly on simple exact-match queries. Development and analysis of schemes allowing more complex query types that are used in practice (e.g. range queries) has remained open.

The work of [25] suggested enabling efficient range queries on encrypted data not by using OPE but so-called *prefix-preserving encryption* (PPE) [33, 5]. Unfortunately, as discussed in [25, 2], PPE schemes are subject to certain attacks in this context; particular queries can completely reveal some of the underlying plaintexts in the database. Moreover, their use necessitates specialized data structures and query formats, which practitioners would prefer to avoid.

Allowing range queries on encrypted data in the public-key setting was studied in [12, 30]. While their schemes provably provide strong security, they are not efficient in our setting, requiring to scan the whole database on every query.

Finally, we clarify that [1], in addition to suggesting the OPE primitive, *does* provide a construction. However, the construction is rather ad-hoc and is designed for a setting where users know all data in advance. Accordingly the encryption algorithm must take as input all the plaintexts in the database. Such setting is not always practical, so a stateless scheme whose encryption algorithm can process single plaintexts on the fly is preferable. Moreover, [1] does not define security nor provide any formal security analysis.

**DEFINING SECURITY OF OPE.** Our first goal is to devise a rigorous definition of security that OPE schemes should satisfy. Of course, such schemes cannot satisfy standard notions of security, such as

indistinguishability against chosen-plaintext attack (IND-CPA), as they are not only deterministic, but also leak the order-relations among the plaintexts. (In particular, an adversary against an OPE scheme that queries two pairs with opposite order can trivially break IND-CPA, as the ciphertexts have the same order as their plaintexts.) So, although we cannot target a notion on the level of IND-CPA, we want to define the best possible security subject to this order-preserving constraint. (Such an approach was taken previously in the case of deterministic public-key encryption [6, 11, 7], on-line ciphers [5], and deterministic authenticated encryption [28].)

**WEAKENING IND-CPA.** One approach is to try to weaken the IND-CPA definition appropriately. Indeed, in the case of deterministic symmetric encryption this was done by [8], which formalizes a notion called *indistinguishability under distinct chosen-plaintext attack* or IND-DCPA. (The notion was subsequently applied to message authentication codes in [4].) Since deterministic encryption leaks equality of plaintexts, IND-DCPA restricts the adversary in the IND-CPA experiment to make queries to its left-right-encryption-oracle of the form  $(x_0^1, x_1^1), \dots, (x_0^q, x_1^q)$  such that  $x_0^1, \dots, x_0^q$  are all distinct and  $x_1^1, \dots, x_1^q$  are all distinct. We generalize this to a notion we call *indistinguishability under ordered chosen-plaintext attack* or IND-OCPA, asking these sequences instead to satisfy the same *order relations*. (See Section 3.2.) Surprisingly, we go on to show that this plausible-looking definition is not useful for us, because it cannot be achieved by an OPE scheme unless the size of its ciphertext space is prohibitively large.

**AN ALTERNATIVE APPROACH.** Instead of trying to further restrict the adversary in the IND-OCPA definition, we turn to an approach along the lines of pseudorandom functions (PRFs) or permutations (PRPs), requiring that no adversary can distinguish between oracle access to the encryption algorithm of the scheme, and a corresponding “ideal” object. In our case the latter is a (uniformly) random order-preserving function on the same domain and range. Since order-preserving functions are injective, it also makes sense to aim for a stronger security notion that additionally gives the adversary oracle access to the decryption algorithm or the inverse function, respectively. We call the resulting notion POPF-CCA for *pseudorandom order-preserving function under chosen-ciphertext attack*.

**TOWARDS A CONSTRUCTION.** After having settled on the POPF-CCA notion, we would naturally like to construct an OPE scheme meeting it. Essentially, the encryption algorithm of such a scheme should behave similarly to an algorithm that samples a random order-preserving function from a specified domain and range on-the-fly (dynamically as new queries are made). (Here we note a connection to implementing huge random objects [19] and lazy-sampling [9].) But it is not immediately clear how this can be done; blockciphers, our usual tool in the symmetric-key setting, do not seem helpful in preserving plaintext order. Our construction takes a different route, borrowing some tools from probability theory. We first uncover a relation between a random order-preserving function and the hypergeometric (HG) and negative hypergeometric (NHG) probability distributions.

**THE CONNECTION TO NHG.** To gain some intuition, first observe that any order-preserving function  $f$  from  $\{1, \dots, M\}$  to  $\{1, \dots, N\}$  can be uniquely represented by a combination of  $M$  out of  $N$  ordered items (see Proposition 4.1). Now let us recall a probability distribution that deals with selections of such combinations. Imagine we have  $N$  balls in a bin, out of which  $M$  are black and  $N - M$  are white. At each step, we draw a ball at random without replacement. Consider the random variable  $Y$  describing the total number of balls in our sample after we collect the  $x$ -th black ball. This random variable follows the so-called negative hypergeometric (NHG) distribution. Us-

ing our representation of an order-preserving function, it is not hard to show that  $f(x)$  for a given point  $x \in \{1, \dots, M\}$  has a NHG distribution over a random choice of  $f$ . Assuming an efficient sampling algorithm for the NHG distribution, this gives a rough idea for a scheme, but there are still many subtleties to take care of.

**HANDLING MULTIPLE POINTS.** First, assigning multiple plaintexts to ciphertexts independently according to the NHG distribution cannot work, because the resulting encryption function is unlikely to even be order-preserving. One could try to fix this by keeping track of all previously encrypted plaintexts and their ciphertexts (in both the encryption and decryption algorithms) and adjusting the parameters of the NHG sampling algorithm appropriately for each new plaintext. But we want a stateless scheme, so it cannot keep track of such previous assignments.

**ELIMINATING THE STATE.** As a first step towards eliminating the state, we show that by assigning ciphertexts to plaintexts in a more organized fashion, the state can actually consist of a static but exponentially long random tape. The idea is that, to encrypt plaintext  $x$ , the encryption algorithm performs a binary search down to  $x$ . That is, it first assigns  $\mathcal{Enc}(K, M/2)$ , then  $\mathcal{Enc}(K, M/4)$  if  $m < M/2$  and  $\mathcal{Enc}(K, 3M/4)$  otherwise, and so on, until  $\mathcal{Enc}(K, x)$  is assigned. Crucially, each ciphertext assignment is made according to the output of the NHG sampling algorithm run on appropriate parameters and *coins from an associated portion of the random tape indexed by the plaintext*. (The decryption algorithm can be defined similarly.) Now, it may not be clear that the resulting scheme induces a *random* order-preserving function from the plaintext to ciphertext space (does its distribution get skewed by the binary search?), but we prove (by induction on the size of the plaintext space) that this is indeed the case.

Of course, instead of making the long random tape the secret key  $K$  for our scheme, we can make it the key for a PRF and generate portions of the tape dynamically as needed. However, coming up with a practical PRF construction to use here requires some care. For efficiency it should be blockcipher-based. Since the size of parameters to the NHG sampling algorithm as well as the number of random coins it needs varies during the binary search, and also because such a construction seems useful in general, it should be both variable input-length (VIL) and variable output-length. Such a construction we call a *length-flexible* (LF)-PRF. We propose a generic construction of an LF-PRF from a VIL-PRF and a (keyless) VOL-PRG (pseudorandom generator). Efficient blockcipher-based VIL-PRFs are known, and we suggest a highly efficient blockcipher-based VOL-PRG that is apparently folklore. POPF-CCA security of the resulting OPE scheme can then be easily proved assuming only standard security (pseudorandomness) of the underlying blockcipher.

**SWITCHING FROM NHG TO HG.** Finally, our scheme needs an efficient sampling algorithm for the NHG distribution. Unfortunately, the existence of such an algorithm seems open. It is known that NHG can be approximated by the negative binomial distribution [27], which in turn can be sampled efficiently [17, 15], and that the approximation improves as  $M$  and  $N$  grow. However, quantifying the quality of approximation for fixed parameters seems difficult.

Instead, we turn to a related probability distribution, namely the hypergeometric (HG) distribution, for which a very efficient exact (not approximated) sampling algorithm is known [23, 24]. In our balls-and-bin model with  $M$  black and  $N - M$  white balls, the random variable  $X$  specifying the number of black balls in our sample as soon as  $y$  balls are picked follows the HG distribution. The scheme based on this distribution, which is the one described in the body of the paper, is rather more involved, but nearly as efficient: instead of  $O(\log M) \cdot T_{\text{NHG}}$  running-time it is  $O(\log N) \cdot T_{\text{HG}}$

(where  $T_{\text{NHG}}, T_{\text{HG}}$  are the running-times of the sampling algorithms for the respective distributions), but we show that it is  $O(\log M) \cdot T_{\text{HG}}$  on average.

We note that the hypergeometric distribution was also used in [20] for sampling pseudorandom permutations and constructing blockciphers for short inputs. The authors of [20] were unaware of the efficient sampling algorithms for HG [23, 24] and provided their own realizations based on general sampling methods.

DISCUSSION. It is important to realize that the “ideal” object in our POPF-CCA definition (a random order-preserving function), and correspondingly our OPE construction meeting it, inherently leak some information about the underlying plaintexts. Characterizing this leakage is an important next step in the study of OPE but is outside the scope of our current paper. (Although we mention that our “big-jump attack” of Theorem 3.1 may provide some insight in this regard.)

The point is that practitioners have indicated their desire to use OPE schemes in order to achieve efficient range queries on encrypted data and are willing to live with its security limitations. In response, we provide a scheme meeting what we believe to be a “best-possible” security notion for OPE. This belief can be justified by noting that it is usually the case that a security notion for a cryptographic object is met by a “random” one (which is sometimes built directly into the definition, as in the case of PRFs and PRPs).

But before one fully understands how not very strong security properties of the ideal object, a random order-preserving function, fit the security needs of applications, we *do not recommend* the practical use of our construction.

ON A MORE GENERAL PRIMITIVE. To allow efficient range queries on encrypted data, it is sufficient to have an order-preserving hash function family  $H$  (not necessarily invertible). The overall OPE scheme would then have secret key  $(K_{\text{Enc}}, K_H)$  where  $K_{\text{Enc}}$  is a key for a normal (randomized) encryption scheme and  $K_H$  is a key for  $H$ , and the encryption of  $x$  would be  $\text{Enc}(K_{\text{Enc}}, x) \parallel H(K_H, x)$  (cf. efficiently searchable encryption (ESE) in [6]). Our security notion (in the CPA case) can also be applied to such  $H$ . In fact, there has been some work on hash functions that are order-preserving or have some related properties [26, 16, 21]. But none of these works are concerned with security in any sense. Since our OPE scheme is efficient and already invertible, we have not tried to build any secure order-preserving hash separately.

ON THE PUBLIC-KEY SETTING. Finally, it is interesting to note that in a public-key setting one cannot expect OPE to provide any privacy at all. Indeed, given a ciphertext  $c$  computed under public key  $pk$ , anyone can decrypt  $c$  via a simple binary-search. In the symmetric-key setting a real-life adversary cannot simply encrypt messages itself, so such an attack is unlikely to be feasible.

## 2 Preliminaries

NOTATION AND CONVENTIONS. We refer to members of  $\{0, 1\}^*$  as strings. If  $x$  is a string then  $|x|$  denotes its length in bits and if  $x, y$  are strings then  $x \parallel y$  denotes an encoding from which  $x, y$  are uniquely recoverable. For  $\ell \in \mathbb{N}$  we denote by  $1^\ell$  the string of  $\ell$  “1” bits. If  $S$  is a set then  $x \stackrel{\$}{\leftarrow} S$  denotes that  $x$  is selected uniformly at random from  $S$ . For convenience, for any  $k \in \mathbb{N}$  we write  $x_1, x_2, \dots, x_k \stackrel{\$}{\leftarrow} S$  as shorthand for  $x_1 \stackrel{\$}{\leftarrow} S, x_2 \stackrel{\$}{\leftarrow} S, \dots, x_n \stackrel{\$}{\leftarrow} S$ . If  $A$  is a randomized algorithm and Coins is the set from where it draws its coins, then we write  $A(x, y, \dots)$  as shorthand for  $R \stackrel{\$}{\leftarrow} \text{Coins}; A(x, y, \dots; R)$ , where the latter denotes the result of running  $A$  on inputs  $x, y, \dots$  and coins  $R$ . And  $a \stackrel{\$}{\leftarrow} A(x, y, \dots)$  means that we assign to  $a$  the output of  $A$  run on inputs  $x, y, \dots$

We denote the probability of event  $A$  by  $\Pr[A]$ . If  $A$  depends on a random variable  $X$ , we write  $\Pr_{X \leftarrow D}[A(X)]$  for the probability of  $A$  for  $X$  sampled randomly from distribution  $D$ . If  $B$  is another event,  $\Pr[A | B]$  denotes the conditional probability of  $A$  given  $B$ . Often, the distribution being used is clear and we omit it, as in  $\Pr_X[A(X)]$  (where  $X \leftarrow D$  is implied). Let  $\mathbb{E}[X]$  denote the expected value of  $X$ . Similarly, we use the notation  $\mathbb{E}_{Y \leftarrow D}[X(Y)]$  or  $\mathbb{E}_Y[X(Y)]$  to emphasize that the expected value is taken over the randomness in selecting related random variable  $Y$  from the distribution  $D$ .

For  $a \in \mathbb{N}$  we denote by  $[a]$  the set  $\{1, \dots, a\}$ . For sets  $X$  and  $Y$ , if  $f: X \rightarrow Y$  is a function, then we call  $X$  the domain,  $Y$  the range, and the set  $\{f(x) \mid x \in X\}$  the image of the function. An adversary is an algorithm. By convention, all algorithms are required to be efficient, meaning run in (expected) polynomial-time in the length of their inputs, and their running-time includes that of any overlying experiment.

**SYMMETRIC ENCRYPTION.** A *symmetric encryption scheme*  $\mathcal{SE} = (\mathcal{K}, \mathcal{Enc}, \mathcal{Dec})$  with associated *plaintext-space*  $\mathcal{D}$  and *ciphertext-space*  $\mathcal{R}$  consists of three algorithms.

- The randomized *key generation algorithm*  $\mathcal{K}$  returns a secret key  $K$ .
- The (possibly randomized) *encryption algorithm*  $\mathcal{Enc}$  takes a secret key  $K$ , descriptions of plaintext and ciphertext-spaces  $\mathcal{D}, \mathcal{R}$  and a plaintext  $m$  to return a ciphertext  $c$ .
- The deterministic *decryption algorithm*  $\mathcal{Dec}$  takes the secret key  $K$ , descriptions of plaintext and ciphertext-spaces  $\mathcal{D}, \mathcal{R}$ , and a ciphertext  $c$  to return a corresponding plaintext  $m$  or a special symbol  $\perp$  indicating that the ciphertext was invalid.

Note that the above syntax differs from the usual one in that we specify the plaintext and ciphertext-spaces  $\mathcal{D}, \mathcal{R}$  explicitly; this is for convenience relative to our specific schemes. We require the usual correctness condition, namely that  $\mathcal{Dec}(K, \mathcal{D}, \mathcal{R}, (\mathcal{Enc}(K, \mathcal{D}, \mathcal{R}, m))) = m$  for all  $K$  output by  $\mathcal{K}$  and all  $m \in \mathcal{D}$ . Finally, we say that  $\mathcal{SE}$  is *deterministic* if  $\mathcal{Enc}$  is deterministic.

**IND-CPA.** Let  $\mathcal{LR}(\cdot, \cdot, b)$  denote the function that on inputs  $m_0, m_1$  returns  $m_b$ . For a symmetric encryption scheme  $\mathcal{SE} = (\mathcal{K}, \mathcal{Enc}, \mathcal{Dec})$  and an adversary  $A$  and  $b \in \{0, 1\}$  consider the following experiment:

$$\begin{aligned} & \mathbf{Experiment} \mathbf{Exp}_{\mathcal{SE}}^{\text{ind-cpa-}b}(A) \\ & K \leftarrow \mathcal{K} \\ & b' \leftarrow A^{\mathcal{Enc}(K, \mathcal{LR}(\cdot, \cdot, b))} \\ & \text{Return } b' . \end{aligned}$$

We require that each query  $(m_0, m_1)$  that  $A$  makes to its oracle satisfies  $|m_0| = |m_1|$ . For an adversary  $A$ , define its *ind-cpa advantage* against  $\mathcal{SE}$  as

$$\mathbf{Adv}_{\mathcal{SE}}^{\text{ind-cpa}}(A) = \Pr \left[ \mathbf{Exp}_{\mathcal{SE}}^{\text{ind-cpa-1}}(A) = 1 \right] - \Pr \left[ \mathbf{Exp}_{\mathcal{SE}}^{\text{ind-cpa-0}}(A) = 1 \right] .$$

We say that  $\mathcal{SE}$  is *indistinguishable under chosen-plaintext attacks* (IND-CPA-secure) if the ind-cpa advantage of any adversary against  $\mathcal{SE}$  is small.

**PSEUDORANDOM FUNCTIONS (PRFs).** We say that  $\mathcal{F} = (\mathcal{K}, F)$  is a *function family* on domain  $\mathcal{D}$  and range  $\mathcal{R}$  if  $\mathcal{K}$  outputs random keys and for each key  $K \leftarrow \mathcal{K}$  the map  $F(K, \cdot)$  is a function from  $\mathcal{D}$  to  $\mathcal{R}$ . We refer to  $F(K, \cdot)$  as an *instance* of  $\mathcal{F}$ .

Let  $\text{Func}_{\mathcal{D},\mathcal{R}}$  denote the set of all functions from  $\mathcal{D}$  to  $\mathcal{R}$ . For any adversary  $A$ , the *prf-advantage* against function family  $\mathcal{F} = (\mathcal{K}, F)$  is defined as

$$\text{Adv}_{\mathcal{F}}^{\text{prf}}(A) = \Pr_{K \xleftarrow{\$} \mathcal{K}} \left[ A^{F(K, \cdot)} = 1 \right] - \Pr_{f \xleftarrow{\$} \text{Func}_{\mathcal{D},\mathcal{R}}} \left[ A^{f(\cdot)} = 1 \right]$$

We say that  $\mathcal{F}$  is a *pseudorandom function* (PRF) if for any efficient adversary  $A$ ,  $\text{Adv}_{\mathcal{F}}^{\text{prf}}(A)$  is small.

### 3 Order-Preserving Encryption and its Security

We begin by defining a primitive for deterministic encryption schemes that preserve order on their plaintext space.

#### 3.1 Order-Preserving Encryption (OPE)

For  $A, B \subseteq \mathbb{N}$  with  $|A| \leq |B|$ , a function  $f: A \rightarrow B$  is *order-preserving* (a.k.a. monotonically increasing) if for all  $i, j \in A$ ,  $f(i) > f(j)$  iff  $i > j$ . We say that deterministic encryption scheme  $\mathcal{SE} = (\mathcal{K}, \mathcal{Enc}, \mathcal{Dec})$  with plaintext and ciphertext spaces  $\mathcal{D}, \mathcal{R}$  is an *order-preserving encryption* (OPE) scheme if  $\mathcal{Enc}(K, \cdot)$  is an order-preserving function from  $\mathcal{D}$  to  $\mathcal{R}$  for all  $K$  output by  $\mathcal{K}$  (with elements of  $\mathcal{D}, \mathcal{R}$  interpreted as numbers, encoded as strings).

For simplicity, we will often assume a plaintext space  $[M]$  and ciphertext space  $[N]$  for some  $N \geq M \in \mathbb{N}$ .

#### 3.2 Security of OPE

OPE obviously cannot be IND-CPA-secure, as it leaks order of plaintexts. A natural question arises: can we weaken the IND-CPA-notion just enough so that it is achievable by an OPE scheme, but is still as strong as possible?

A FIRST TRY. Security of deterministic symmetric encryption was introduced in [8], as a notion they call *security under distinct chosen-plaintext attack (IND-DCPA)*. (It will not be important to consider CCA now.) The idea is that because deterministic encryption leaks plaintext equality, the adversary  $A$  in the IND-CPA experiment defined in Section 2 is restricted to make only *distinct* queries on either side of its oracle (as otherwise there is a trivial attack). That is, supposing  $A$  makes queries  $(m_0^1, m_1^1), \dots, (m_0^q, m_1^q)$ , they require that  $m_b^1, \dots, m_b^q$  are all distinct for  $b \in \{0, 1\}$ .

Noting that any OPE scheme analogously leaks order relations of plaintexts, consider extending the above approach to take this into account. In particular, let us further require the above queries made by  $A$  to have the same “order pattern.”

IND-OCPA. For a symmetric order-preserving encryption scheme  $\text{OPE} = (\mathcal{K}, \mathcal{Enc}, \mathcal{Dec})$  and an adversary  $A$  and  $b \in \{0, 1\}$  consider the following experiment:

**Experiment**  $\text{Exp}_{\text{OPE}}^{\text{ind-ocpa-}b}(A)$   
 $K \xleftarrow{\$} \mathcal{K}$   
 $d \xleftarrow{\$} A^{\mathcal{Enc}(K, \mathcal{LR}(\cdot, b))}$   
 Return  $d$



We require that each query  $(m_0, m_1)$  that  $A$  makes to its oracle satisfies  $|m_0| = |m_1|$ , and also that the LR-queries have the same *order pattern*, i.e.  $m_0^i < m_0^j$  iff  $m_1^i < m_1^j$  for all  $1 \leq i, j \leq q$ . For an adversary  $A$ , define its *indistinguishability under ordered chosen-plaintext attack (IND-OCPA) advantage* against OPE as

$$\mathbf{Adv}_{\text{OPE}}^{\text{ind-ocpa}}(A) = \Pr \left[ \mathbf{Exp}_{\text{OPE}}^{\text{ind-ocpa-1}}(A) = 1 \right] - \Pr \left[ \mathbf{Exp}_{\text{OPE}}^{\text{ind-ocpa-0}}(A) = 1 \right].$$

IND-OCPA IS NOT USEFUL. Defining IND-OCPA adversary seems like a plausible way to analyze security for OPE. Surprisingly, it turns out to be not useful for us. In the following theorem, we show that IND-OCPA is unachievable by a practical order-preserving encryption scheme, in that an OPE scheme cannot be IND-OCPA unless its ciphertext space is extremely large, namely, super-exponential in the size of the plaintext space. [Note: this is an improvement over the corresponding result in the proceedings version [10] of this paper, in that it holds also for ciphertext spaces that have size exponential in the plaintext space size.]

**Theorem 3.1.** *Let  $\text{OPE} = (\mathcal{K}, \mathcal{Enc}, \mathcal{Dec})$  be an order-preserving encryption scheme on plaintext-space  $[M]$  and ciphertext-space  $[N]$ , where  $N < t^{\lfloor M/4 \rfloor}$  for some constant integer  $t > 1$ . There exists an IND-OCPA adversary  $A$  against OPE such that*

$$\mathbf{Adv}_{\text{OPE}}^{\text{ind-ocpa}}(A) > \frac{1}{2t}.$$

Furthermore,  $A$  runs in time  $O(\log N)$  and makes at most 3 oracle queries.

*Proof.* Let  $M' = \lfloor M/2 \rfloor$ . Consider the following IND-OCPA adversary  $A$  against OPE:

**Adversary**  $A^{\mathcal{Enc}(K, \mathcal{LR}(\cdot, b))}$   
 $m_0 \xleftarrow{\$} [M'], m_1 \leftarrow M - m_0 + 1$   
 $c \leftarrow \mathcal{Enc}(K, \mathcal{LR}(m_0, m_1, b))$   
 $c_L \leftarrow 0 ; c_R \leftarrow N + 1$   
 If  $m_0 > 1$ :  
 $m_L \leftarrow m_0 - 1, m_R \leftarrow m_1 + 1$   
 $c_L \leftarrow \mathcal{Enc}(K, \mathcal{LR}(m_L, m_L, b))$   
 $c_R \leftarrow \mathcal{Enc}(K, \mathcal{LR}(m_R, m_R, b))$   
 Return 1 with probability  $\frac{c - c_L}{c_R - c_L}$   
 Else return 0

The IND-OCPA correctness and efficiency claims of  $A$  should be clear from the construction.

Fix a key  $K$ , so that  $\mathcal{Enc}(K, \cdot)$  is a well-defined order-preserving function from  $[M]$  to  $[N]$ . For  $m \in [M']$ , let  $X_m = \mathcal{Enc}(K, M - m + 1) - \mathcal{Enc}(K, m)$  and  $X_0 = N + 1$ . Let  $S$  be the set of messages  $m$  in  $[M']$  such that  $\frac{X_m}{X_{m-1}} \leq \frac{1}{t}$ . Then if  $|S| > M'/2$  we have

$$N + 1 = X_0 \geq \frac{X_0}{X_{M'}} = \prod_{m \in [M']} \frac{X_{m-1}}{X_m} > t^{M'/2} \geq t^{\lfloor M/4 \rfloor},$$

a contradiction to  $N < t^{\lfloor M/4 \rfloor}$ . Thus,  $|S| \leq M'/2$  and so

$$\Pr \left[ m \xleftarrow{\$} [M'] \mid \frac{X_m}{X_{m-1}} \leq \frac{1}{t} \right] \leq \frac{1}{2}.$$

In the following, for given  $m_0 \in [M']$  let  $m_1, m_L, m_R, c_0, c_1, c_L, c_R$  be determined from  $m_0$  as in A. We have

$$\begin{aligned}
\mathbf{Adv}_{\text{OPE}}^{\text{ind-ocpa}}(A) &= \Pr \left[ \mathbf{Exp}_{\text{OPE}}^{\text{ind-ocpa-1}}(A) = 1 \right] - \Pr \left[ \mathbf{Exp}_{\text{OPE}}^{\text{ind-ocpa-0}}(A) = 1 \right] \\
&= \Pr \left[ A^{\mathcal{Enc}(K, \mathcal{LR}(\cdot, 1))} = 1 \right] - \Pr \left[ A^{\mathcal{Enc}(K, \mathcal{LR}(\cdot, 0))} = 1 \right] \\
&= \mathbb{E}_{m_0 \xleftarrow{\$} [M']} \left[ \frac{c_1 - c_L}{c_R - c_L} \right] - \mathbb{E}_{m_0 \xleftarrow{\$} [M']} \left[ \frac{c_0 - c_L}{c_R - c_L} \right] \\
&= \mathbb{E}_{m_0 \xleftarrow{\$} [M']} \left[ \frac{c_1 - c_0}{c_R - c_L} \right] \\
&\geq \frac{1}{t} \Pr_{m_0 \xleftarrow{\$} [M']} \left[ \frac{c_1 - c_0}{c_R - c_L} \geq \frac{1}{t} \right] \\
&= \frac{1}{t} \Pr_{m_0 \xleftarrow{\$} [M']} \left[ \frac{X_m}{X_{m-1}} \geq \frac{1}{t} \right] \\
&> \frac{1}{2t}.
\end{aligned}$$

This completes the proof.  $\square$

DISCUSSION. Obviously, to have ciphertext space size super-exponential in that of the plaintext space would be inconceivable, so for all intents and purposes, Theorem 3.1 shows IND-OCPA is unachievable for all practical OPE schemes.

The adversary in the proof of Theorem 3.1 uses what we call a “big-jump attack” to distinguish between ciphertexts of messages that are “very close” and “far apart.” The attack shows that *any* practical OPE scheme inherently leaks more information about the plaintexts than just their ordering, namely some information about their relative distances. We return to this point later.

AN ALTERNATIVE APPROACH. Since OPE inherently leaks distance information about plaintexts, further weakening of IND-CPA does not seem very fruitful, as long as attacks can still sample far-apart versus close-together plaintexts. Instead, we take the approach used in defining security e.g. of pseudorandom permutations (PRPs) [18] or on-line PRPs [5], where one asks that oracle access to the function in question be indistinguishable from access to the corresponding “ideal” random object, e.g. a random permutation or a random on-line permutation. As order-preserving functions are injective, we consider the “strong” version of such a definition where an inverse oracle is also given.

POPF-CCA. Fix an order-preserving encryption scheme  $\mathcal{SE} = (\mathcal{K}, \mathcal{Enc}, \mathcal{Dec})$  with plaintext space  $\mathcal{D}$  and ciphertext space  $\mathcal{R}$ ,  $|\mathcal{D}| \leq |\mathcal{R}|$ . For an adversary  $A$  against  $\mathcal{SE}$ , define its *pseudorandom order-preserving function advantage under chosen-ciphertext attacks (POPF-CCA) advantage* against  $\mathcal{SE}$  as

$$\mathbf{Adv}_{\mathcal{SE}}^{\text{popf-cca}}(A) = \Pr_{K \xleftarrow{\$} \mathcal{K}} \left[ A^{\mathcal{Enc}(K, \cdot), \mathcal{Dec}(K, \cdot)} = 1 \right] - \Pr_{g \xleftarrow{\$} \text{OPF}_{\mathcal{D}, \mathcal{R}}} \left[ A^{g(\cdot), g^{-1}(\cdot)} = 1 \right],$$

where  $\text{OPF}_{\mathcal{D}, \mathcal{R}}$  denotes the set of all order-preserving functions from  $\mathcal{D}$  to  $\mathcal{R}$ . We say  $\mathcal{SE}$  is POPF-secure if this advantage is small.

LAZY SAMPLING. Now in order for this notion to be useful, i.e. to be able to show that a scheme achieves it, we also need a way to implement  $A$ ’s oracles in the “ideal” experiment efficiently. In

other words, we need to show how to “lazy sample” (a term from [9]) a random order-preserving function and its inverse.<sup>1</sup>

As shown in [9], lazy sampling of “exotic” functions with many constraints can be tricky. In the case of a random order-preserving function, it turns out that straightforward procedures—which assign a random point in the range to a queried domain point, subject to the obvious remaining constraints—do not work (that is, the resulting function is not uniformly distributed over the set of all such functions). So how can we lazy sample such a function, if it is possible at all? We address this issue next.

A CAVEAT. Before proceeding, we note that a shortcoming of our POPF-CCA notion is it does not lead to a nice answer to the question of what information about the data is leaked by a secure OPE scheme, but only reduces this to the question of what information the “ideal object” (a random order-preserving function) leaks. Although practitioners have indicated that they are willing to live with the security limitations of OPE for its useful functionality, more precisely characterizing the latter remains an important next step before our schemes should be considered for practical deployment.

## 4 Lazy Sampling a Random Order-Preserving Function

In this section, we show how to lazy-sample a random order-preserving function and its inverse. This result may also be of independent interest, since the more general question of what functions can be lazy-sampled is interesting in its own right, and it may find other applications as well, e.g. to [29]. We first uncover a connection between a random order-preserving function and the hypergeometric (HG) probability distribution.

### 4.1 The Hypergeometric Connection

To gain some intuition we start with the following claim.

**Proposition 4.1.** *There is bijection between the set  $\text{OPF}_{\mathcal{D},\mathcal{R}}$  containing all order-preserving functions from a domain  $\mathcal{D}$  of size  $M$  to a range  $\mathcal{R}$  of size  $N \geq M$  and the set of all possible combinations of  $M$  out of  $N$  ordered items.*

*Proof.* Without loss of generality, it is enough to prove the result for domain  $[M]$  and range  $[N]$ . Imagine a graph with its  $x$ -axis marked with integers from 1 to  $M$  and its  $y = f(x)$ -axis marked with integers from 1 to  $N$ . Given  $S$ , a set of  $M$  distinct integers from  $[N]$ , construct an order-preserving function from  $[M]$  to  $[N]$  by mapping each  $i \in [M]$  to the  $i$ th smallest element in  $S$ . So, an  $M$ -out-of- $N$  combination corresponds to a unique order-preserving function. On the other hand, consider an order-preserving function  $f$  from  $[M]$  to  $[N]$ . The image of  $f$  defines a set of  $M$  distinct objects in  $[N]$ , so an order-preserving function corresponds to a unique  $M$ -out-of- $N$  combination.  $\square$

---

<sup>1</sup>For example, in the case of a random function from the set of *all* functions one can simply assign a random point from the range to each new point queried from the domain. In the case of a random permutation, the former can be chosen from the set of all previously unassigned points in the range, and lazy sampling of its inverse can be done similarly. A lazy sampling procedure for a random on-line PRP and its inverse via a tree-based characterization was given in [5].

Using the above combination-based characterization it is straightforward to justify the following equality, defined for  $M, N \in \mathbb{N}$  and any  $x, x + 1 \in [M], y \in [N]$ :

$$\Pr_{f \stackrel{s}{\leftarrow} \text{OPF}_{[M],[N]}} [f(x) \leq y < f(x + 1)] = \frac{\binom{y}{x} \cdot \binom{N-y}{M-x}}{\binom{N}{M}}. \quad (1)$$

Now let us recall a particular distribution dealing with an experiment of selecting from combinations of items.

**HYPERGEOMETRIC DISTRIBUTION.** Consider the following balls-and-bins model. Assume we have  $N$  balls in a bin out of which  $M$  balls are black and  $N - M$  balls are white. At each step we draw a ball at random, without replacement. Consider a random variable  $X$  that describes the number of black balls chosen after a *sample size* of  $y$  balls are picked. This random variable has a hypergeometric distribution, and the probability that  $X = x$  for the parameters  $N, M, y$  is

$$P_{\text{HG}}(x; N, M, y) = \frac{\binom{y}{x} \cdot \binom{N-y}{M-x}}{\binom{N}{M}}. \quad (2)$$

Intuitively, Equations 1 and 2 imply that we can construct a random order-preserving function  $f$  from  $[M]$  to  $[N]$  as an experiment involving  $N$  balls,  $M$  of which are black. Choosing balls randomly without replacement, if the  $y$ -th ball we pick is black then the least unmapped point in the domain is mapped to  $y$  under  $f$ . Of course, this experiment is too inefficient to be performed directly. But we will use the hypergeometric distribution to design procedures that efficiently and recursively lazy sample a random order-preserving function and its inverse.

## 4.2 The LazySample Algorithms

Here we give our algorithms **LazySample**, **LazySampleInv** that lazy sample a random order-preserving function from domain  $\mathcal{D}$  to range  $\mathcal{R}$ ,  $|\mathcal{D}| \leq |\mathcal{R}|$ , and its inverse, respectively. The algorithms share and maintain joint state. We assume that both  $\mathcal{D}$  and  $\mathcal{R}$  are sets of consecutive integers.

**TWO SUBROUTINES.** Our algorithms make use of two subroutines. The first, denoted **HG**, takes inputs  $M, N$ , and  $y \in \{0, 1, \dots, N\}$  to return  $x \in \{0, 1, \dots, M\}$  such that for each  $x^* \in \{0, 1, \dots, M\}$  we have  $x = x^*$  with probability  $P_{\text{HG}}(x; N, M, y)$  over the coins of **HG**. (Efficient algorithms for this exist, and we discuss them in Section 4.5.) The second, denoted **GetCoins**, takes inputs  $1^\ell, \mathcal{D}, \mathcal{R}$ , and  $b||z$ , where  $b \in \{0, 1\}$  and  $z \in \mathcal{R}$  if  $b = 0$  and  $z \in \mathcal{D}$  otherwise, to return  $cc \in \{0, 1\}^\ell$ .

**THE ALGORITHMS.** To define our algorithms, let us denote by  $w \stackrel{cc}{\leftarrow} S$  that  $w$  is assigned a value sampled uniformly at random from set  $S$  using coins  $cc$  of length  $\ell_S$ , where  $\ell_S$  denotes the number of coins needed to do so. Let  $\ell_1 = \ell(M, N, y - r)$  denote the number of coins needed by **HG** on inputs  $M, N, y - r$ . Our algorithms are given in Figure 1. Note that the arrays  $F, I$ , initially empty, are global and shared between the algorithms; also, for now, think of **GetCoins** as returning fresh random coins. We later implement it by using a PRF on the same parameters to eliminate the joint state.

**OVERVIEW.** To determine the image of input  $m$ , **LazySample** employs a strategy of mapping “range gaps” to “domain gaps” in a recursive, binary search manner. By “range gap” or “domain

<p><b>LazySample</b>(<math>\mathcal{D}, \mathcal{R}, m</math>)</p> 01 $M \leftarrow  \mathcal{D}  ; N \leftarrow  \mathcal{R} $ 02 $d \leftarrow \min(\mathcal{D}) - 1 ; r \leftarrow \min(\mathcal{R}) - 1$ 03 $y \leftarrow r + \lceil N/2 \rceil$ 04 If $ \mathcal{D}  = 1$ then 05     If $F[\mathcal{D}, \mathcal{R}, m]$ is undefined then 06 $cc \stackrel{\$}{\leftarrow} \text{GetCoins}(1^{\ell_{\mathcal{R}}}, \mathcal{D}, \mathcal{R}, 1  m)$ 07 $F[\mathcal{D}, \mathcal{R}, m] \stackrel{cc}{\leftarrow} \mathcal{R}$ 08     Return $F[\mathcal{D}, \mathcal{R}, m]$ 09 If $I[\mathcal{D}, \mathcal{R}, y]$ is undefined then 10 $cc \stackrel{\$}{\leftarrow} \text{GetCoins}(1^{\ell_1}, \mathcal{D}, \mathcal{R}, 0  y)$ 11 $I[\mathcal{D}, \mathcal{R}, y] \stackrel{\$}{\leftarrow} \text{HG}(M, N, y - r; cc)$ 12 $x \leftarrow d + I[\mathcal{D}, \mathcal{R}, y]$ 13 If $m \leq x$ then 14 $\mathcal{D} \leftarrow \{d + 1, \dots, x\}$ 15 $\mathcal{R} \leftarrow \{r + 1, \dots, y\}$ 16 Else 17 $\mathcal{D} \leftarrow \{x + 1, \dots, d + M\}$ 18 $\mathcal{R} \leftarrow \{y + 1, \dots, r + N\}$ 19 Return <b>LazySample</b> ( $\mathcal{D}, \mathcal{R}, m$ )	<p><b>LazySampleInv</b>(<math>\mathcal{D}, \mathcal{R}, c</math>)</p> 20 $M \leftarrow  \mathcal{D}  ; N \leftarrow  \mathcal{R} $ 21 $d \leftarrow \min(\mathcal{D}) - 1 ; r \leftarrow \min(\mathcal{R}) - 1$ 22 $y \leftarrow r + \lceil N/2 \rceil$ 23 If $ \mathcal{D}  = 1$ then $m \leftarrow \min(\mathcal{D})$ 24     If $F[\mathcal{D}, \mathcal{R}, m]$ is undefined then 25 $cc \stackrel{\$}{\leftarrow} \text{GetCoins}(1^{\ell_{\mathcal{R}}}, \mathcal{D}, \mathcal{R}, 1  m)$ 26 $F[\mathcal{D}, \mathcal{R}, m] \stackrel{cc}{\leftarrow} \mathcal{R}$ 27     If $F[\mathcal{D}, \mathcal{R}, m] = c$ then return $m$ 28     Else return $\perp$ 29 If $I[\mathcal{D}, \mathcal{R}, y]$ is undefined then 30 $cc \stackrel{\$}{\leftarrow} \text{GetCoins}(1^{\ell_1}, \mathcal{D}, \mathcal{R}, 0  y)$ 31 $I[\mathcal{D}, \mathcal{R}, y] \stackrel{\$}{\leftarrow} \text{HG}(M, N, y - r; cc)$ 32 $x \leftarrow d + I[\mathcal{D}, \mathcal{R}, y]$ 33 If $c \leq y$ then 34 $\mathcal{D} \leftarrow \{d + 1, \dots, x\}$ 35 $\mathcal{R} \leftarrow \{r + 1, \dots, y\}$ 36 Else 37 $\mathcal{D} \leftarrow \{x + 1, \dots, d + M\}$ 38 $\mathcal{R} \leftarrow \{y + 1, \dots, r + N\}$ 39 Return <b>LazySampleInv</b> ( $\mathcal{D}, \mathcal{R}, c$ )
--	---

Figure 1: Algorithms **LazySample** and **LazySampleInv** for lazy-sampling a pseudorandom order-preserving function and its inverse by sampling the hypergeometric distribution.

gap,” we mean an imaginary barrier between two consecutive points in the range or domain, respectively. When run, the algorithm first maps the middle range gap  $y$  (the gap between the middle two range points) to a domain gap. To determine the mapping, on line 11 it sets, according to the hypergeometric distribution, how many points in  $\mathcal{D}$  are mapped up to range point  $y$  and stores this value in array  $I$ . (In the future the array is referenced instead of choosing this value anew.) Thus we have that  $f(x) \leq y < f(x + 1)$  (cf. (1)), where  $x = d + I[\mathcal{D}, \mathcal{R}, y]$  as computed on line 12. So, we can view the range gap between  $y$  and  $y + 1$  as having been mapped to the domain gap between  $x$  and  $x + 1$ .

If the input domain point  $m$  is below (resp. above) the domain gap, the algorithm recurses on line 19 on the lower (resp. upper) half of the range and the lower (resp. upper) part of the domain, mapping further “middle” range gaps to domain gaps. This process continues until the gaps on either side of  $m$  have been mapped to by some range gaps. Finally, on line 07, the algorithm samples a range point uniformly at random from the “window” defined by the range gaps corresponding to  $m$ ’s neighboring domain gaps. This result is assigned to array  $F$  as the image of  $m$  under the lazy-sampled function.

### 4.3 Correctness

When `GetCoins` returns truly random coins, it is clear that **LazySample** and **LazySampleInv** are consistent and sample an order-preserving function and its inverse respectively. But we need a stronger claim, namely, that our algorithms sample a *random* order-preserving function and its

inverse. We show this by arguing that any (even computationally unbounded) adversary has no advantage in distinguishing oracle access to a random order-preserving function and its inverse from that to the algorithms **LazySample**, **LazySampleInv**. The following theorem states this claim.

**Theorem 4.2.** *Suppose `GetCoins` returns truly random coins on each new input. Then for any (even computationally unbounded) algorithm  $A$  we have*

$$\Pr_{g \xleftarrow{\$} \text{OPF}_{\mathcal{D}, \mathcal{R}}} \left[ A^{g(\cdot), g^{-1}(\cdot)} = 1 \right] = \Pr \left[ A^{\mathbf{LazySample}(\mathcal{D}, \mathcal{R}, \cdot), \mathbf{LazySampleInv}(\mathcal{D}, \mathcal{R}, \cdot)} = 1 \right],$$

where  $g^{-1}$  denotes the inverse of OPF  $g$ .

*Proof.* Since we consider unbounded adversaries, we can ignore the inverse oracle in our analysis, since such an adversary can always query all points in the domain to learn all points in the image. Let  $M = |\mathcal{D}|$ ,  $N = |\mathcal{R}|$ ,  $d = \min(\mathcal{D}) - 1$ , and  $r = \min(\mathcal{R}) - 1$ . We will say that two functions  $g, h : \mathcal{D} \rightarrow \mathcal{R}$  are *equivalent* if  $g(m) = h(m)$  for all  $m \in \mathcal{D}$ . (Note that if  $\mathcal{D} = \emptyset$ , any two functions  $g, h : \mathcal{D} \rightarrow \mathcal{R}$  are vacuously equivalent.) Let  $f$  be any function in  $\text{OPF}_{\mathcal{D}, \mathcal{R}}$ . To prove the theorem, it is enough to show that the function defined by **LazySample**  $(\mathcal{D}, \mathcal{R}, \cdot)$  is equivalent to  $f$  with probability  $1/|\text{OPF}_{\mathcal{D}, \mathcal{R}}|$ . We prove this using strong induction on  $M$  and  $N$ .

Consider the base case where  $M = 1$ , i.e.,  $\mathcal{D} = \{m\}$  for some  $m$ , and  $N \geq M$ . When it is first called, **LazySample**  $(\mathcal{D}, \mathcal{R}, m)$  will determine an element  $c$  uniformly at random from  $\mathcal{R}$  and enter it into  $F[\mathcal{D}, \mathcal{R}, m]$ , whereupon any future calls of **LazySample**  $(\mathcal{D}, \mathcal{R}, m)$  will always output  $F[\mathcal{D}, \mathcal{R}, m] = c$ . Thus, the output of **LazySample**  $(\mathcal{D}, \mathcal{R}, m)$  is always  $c$ , so **LazySample**  $(\mathcal{D}, \mathcal{R}, \cdot)$  is equivalent to  $f$  if and only if  $c = f(m)$ . Since  $c$  is chosen randomly from  $\mathcal{R}$ ,  $c = f(m)$  with probability  $1/|\mathcal{R}|$ . Thus, **LazySample**  $(\mathcal{D}, \mathcal{R}, m)$  is equivalent to  $f(m)$  with probability  $1/|\mathcal{R}| = 1/|\text{OPF}_{\mathcal{D}, \mathcal{R}}|$ .

Now suppose  $M > 1$ , and  $N \geq M$ . As an induction hypothesis, assume that for all domains  $\mathcal{D}'$  of size  $M'$  and ranges  $\mathcal{R}'$  of size  $N' \geq M'$ , where either  $M' < M$  or  $(M' = M \text{ and } N' < N)$ , and for any function  $f'$  in  $\text{OPF}_{\mathcal{D}', \mathcal{R}'}$ , **LazySample**  $(\mathcal{D}', \mathcal{R}', \cdot)$  is equivalent to  $f'$  with probability  $1/|\text{OPF}_{\mathcal{D}', \mathcal{R}'}|$ .

When it is first called, **LazySample**  $(\mathcal{D}, \mathcal{R}, \cdot)$  sets  $I[\mathcal{D}, \mathcal{R}, y]$  to be the value of  $\text{HG}(M, N, y - r; cc)$ , where  $y = r + \lceil N/2 \rceil$ ,  $r = \min(\mathcal{R}) - 1$ . Henceforth, on this and future calls of **LazySample**  $(\mathcal{D}, \mathcal{R}, m)$ , the algorithm sets  $x = d + I[\mathcal{D}, \mathcal{R}, y - r]$  and runs **LazySample**  $(\mathcal{D}_1, \mathcal{R}_1, m)$  if  $m \leq x$ , or run **LazySample**  $(\mathcal{D}_2, \mathcal{R}_2, m)$  if  $m > x$ , where  $\mathcal{D}_1 = \{1, \dots, x\}$ ,  $\mathcal{R}_1 = \{1, \dots, y\}$ ,  $\mathcal{D}_2 = \{x + 1, \dots, M\}$ ,  $\mathcal{R}_2 = \{y + 1, \dots, N\}$ . Let  $f_1$  be  $f$  restricted to the domain  $\mathcal{D}_1$ , and let  $f_2$  be  $f$  restricted to the domain  $\mathcal{D}_2$ . Let  $x_0$  be the unique integer in  $\mathcal{D} \cup \{d\}$  such that  $f(z) \leq y$  for all  $z \in \mathcal{D}$  with  $z \leq x_0$ , and  $f(z) > y$  for all  $z \in \mathcal{D}$  with  $z > x_0$ . Note then that **LazySample**  $(\mathcal{D}, \mathcal{R}, \cdot)$  is equivalent to  $f$  if and only if all three of the following events occur:

$E_1$ :  $f$  restricted to range  $\mathcal{R}_1$  stays within domain  $\mathcal{D}_1$ , and  $f$  restricted to range  $\mathcal{R}_2$  stays within domain  $\mathcal{D}_2$ —that is,  $x$  is chosen to be  $x_0$ .

$E_2$ : **LazySample**  $(\mathcal{D}_1, \mathcal{R}_1, \cdot)$  is equivalent to  $f_1$ .

$E_3$ : **LazySample**  $(\mathcal{D}_2, \mathcal{R}_2, \cdot)$  is equivalent to  $f_2$ .

By the law of conditional probability, and since  $E_2$  and  $E_3$  are independent,

$$\begin{aligned} \Pr[E_1 \cap E_2 \cap E_3] &= \Pr[E_1] \Pr[E_2 \cap E_3 \mid E_1] \\ &= \Pr[E_1] \Pr[E_2 \mid E_1] \Pr[E_3 \mid E_1]. \end{aligned}$$

$\Pr[E_1]$  is the hypergeometric probability that  $\text{HG}(M, N, y - r)$  will return  $x_0 - d$ , so

$$\Pr[E_1] = P_{\text{HG}}(x_0 - d; N, M, \lceil N/2 \rceil) = \frac{\binom{\lceil N/2 \rceil}{x_0 - d} \binom{N - \lceil N/2 \rceil}{M - (x_0 - d)}}{\binom{N}{M}}.$$

Assuming for the moment that neither  $\mathcal{D}_1$  nor  $\mathcal{D}_2$  are empty, notice that both  $|\mathcal{R}_1|$  and  $|\mathcal{R}_2|$  are strictly less than  $|\mathcal{R}|$ , and  $|\mathcal{D}_1|$  and  $|\mathcal{D}_2|$  are less than or equal to  $|\mathcal{D}|$ , so the induction hypothesis holds for each. That is, **LazySample** ( $\mathcal{D}_1, \mathcal{R}_1, \cdot$ ) is equivalent to  $f_1$  with probability  $1/|\text{OPF}_{\mathcal{D}_1, \mathcal{R}_1}| = 1/\binom{|\mathcal{R}_1|}{|\mathcal{D}_1|}$ , and **LazySample** ( $\mathcal{D}_2, \mathcal{R}_2, \cdot$ ) is equivalent to  $f_2$  with probability  $1/|\text{OPF}_{\mathcal{D}_2, \mathcal{R}_2}| = 1/\binom{|\mathcal{R}_2|}{|\mathcal{D}_2|}$ . Thus, we have that

$$\Pr[E_2 | E_1] = \frac{1}{\binom{\lceil N/2 \rceil}{x_0 - d}} \quad \text{and} \quad \Pr[E_3 | E_1] = \frac{1}{\binom{N - \lceil N/2 \rceil}{d + M - x_0}}.$$

Also, note that if  $\mathcal{D}_1 = \emptyset$ , then  $\Pr[E_2 | E_1] = 1 = \frac{1}{\binom{\lceil N/2 \rceil}{x_0 - d}}$  since  $x_0 = d$ . Likewise, if  $\mathcal{D}_2 = \emptyset$ , then  $\Pr[E_3 | E_1]$  will be the same as above. We conclude that

$$\Pr[E_1 \cap E_2 \cap E_3] = \frac{\binom{\lceil N/2 \rceil}{x_0 - d} \binom{N - \lceil N/2 \rceil}{M - (x_0 - d)}}{\binom{N}{M}} \cdot \frac{1}{\binom{\lceil N/2 \rceil}{x_0 - d}} \cdot \frac{1}{\binom{N - \lceil N/2 \rceil}{d + M - x_0}} = \frac{1}{\binom{N}{M}}.$$

Thus, **LazySample** ( $\mathcal{D}, \mathcal{R}, \cdot$ ) is equivalent to  $f$  with probability  $\frac{1}{\binom{N}{M}} = \frac{1}{|\text{OPF}_{\mathcal{D}, \mathcal{R}}|}$ . Since  $f$  was an arbitrary element of  $\text{OPF}_{\mathcal{D}, \mathcal{R}}$ , the result follows.  $\square$

We clarify that in the theorem,  $A$ 's oracles for **LazySample**, **LazySampleInv** in the right-hand-side experiment share and update joint state. It is straightforward to check, via simple probability calculations, that the theorem holds for an adversary  $A$  that makes one query. The case of multiple queries is harder. The reason is that the distribution of the responses given to subsequent queries depends on which queries  $A$  has already made, and this distribution is difficult to compute directly. Instead our proof uses strong induction in a way that parallels the recursive nature of our algorithms.

#### 4.4 Efficiency

We characterize efficiency of our algorithms in terms of the number of recursive calls made by **LazySample** or **LazySampleInv** before termination. (The proposition below is just stated in terms of **LazySample** for simplicity; the analogous result holds for **LazySampleInv**.)

**Proposition 4.3.** *The number of recursive calls made by **LazySample** is at most  $\log N + 1$  in the worst-case and at most  $5 \log M + 12$  on average.*

*Proof.* For the worst case bound, note that **LazySample** performs a binary search over the range to map the input domain point, on each recursion cutting the size of the possible range in half. Note that, by the nature of the hypergeometric probabilities, the size of the domain in each iteration can never exceed the size of the range. Thus, when the algorithm is called on a range of size 1, its domain is also of size 1, and the algorithm must terminate. Over the course of  $\log N$  binary-search

recursions, the range will shrink to size 1, so we conclude that a worst-case  $\log N + 1$  recursions are required for **LazySample** to terminate.

For the average case bound, we use a result of Chvátal [13] that the tail of the hypergeometric distribution can be bounded so that

$$\sum_{i=k+1}^M P_{\text{HG}}(i; N, M, c) \leq e^{-2t^2 M},$$

where  $t$  is a fraction such that  $0 \leq t \leq 1 - c/N$ , and  $k = (c/N + t)M$ . Taking  $c = N/2$ , this implies an upper bound on the probability of the hypergeometric distribution assigning our middle domain gap to an “outlying” domain gap:

$$\sum_{i \notin S} P_{\text{HG}}(i; N, M, N/2) \leq 2e^{-2t^2 M} \tag{3}$$

where  $S$  is the subdomain  $[(1/2 - t)M, (1/2 + t)M]$ .

For  $M < 12$ , after at most 12 calls to **LazySample** we will reach a domain of size 1, and terminate. So suppose that  $M \geq 12$ . Taking  $t = 1/4$  in (3) implies that **LazySample** assigns the middle ciphertext gap to a plaintext gap in the “middle subdomain”  $[M/4, 3M/4]$  with probability at least  $1 - 2e^{-2(1/4)^2 M} \geq 1 - 2e^{-3/2} > 1/2$ . When a domain gap in  $S$  is chosen it shrinks the current domain by a fraction of at least  $3/4$ . So, picking in the middle subdomain  $\log_{4/3} M = \frac{\log M}{\log 4/3} < 2.5 \log M$  times will shrink it to size less than 12. Since the probability to pick in the middle subdomain is greater than  $1/2$  on each recursive call of **LazySample**, we expect at most  $5 \log M$  recursive calls to reach domain size  $M < 12$ . Therefore, in total at most  $5 \log M + 12$  recursive calls are needed on average to map an input domain point.  $\square$

Note that the algorithms make one call to HG on each recursion, so an upper-bound on their running-times is then at most  $(\log N + 1) \cdot T_{\text{HG}}$  in the worst-case and at most  $(5 \log M + 12) \cdot T_{\text{HG}}$  on average, where  $T_{\text{HG}}$  denotes the running-time of HG on inputs of size at most  $\log N$ . However, this does not take into account the fact that the size of these inputs decrease on each recursion. Thus, better bounds may be obtained by analyzing the running-time of a specific realization of HG.

## 4.5 Realizing HG

Kachitvichyanukul and Schmeiser [23] designed an efficient implementation of a sampling algorithm HG for the hypergeometric distribution. Their algorithm is exact; it is not an approximation by a related distribution. It is implemented in Wolfram Mathematica and other libraries, and is fast even for large parameters. However, on small parameters the algorithms of [31] perform better. Since the parameter size to HG in our **LazySample** algorithms shrinks across the recursive calls from large to small, it could be advantageous to switch algorithms at some threshold. We refer the reader to [31, 23, 24, 15] for more details.

We comment that the algorithms of [23] are technically only “exact” when the underlying floating-point operations can be performed to infinite precision. In practice, one has to be careful of truncation error. For simplicity, Theorem 4.2 does not take this into account, as in theory the error can be made arbitrarily small by increasing the precision of floating-point operations (independently of  $M, N$ ). But we make this point explicit in Theorem 5.3 where we analyze security of our actual scheme.



## 5 Our OPE Scheme and its Analysis

Algorithms **LazySample**, **LazySampleInv** cannot be directly converted into encryption and decryption procedures because they share and update a joint state, namely arrays  $F$  and  $I$ , which store the outputs of the randomized algorithm **HG**. For our actual scheme, we can eliminate this shared state by implementing the subroutine **GetCoins** (which produces coins for **HG**) as a PRF, and re-constructing entries of  $F$  and  $I$  on-the-fly as needed. However, coming up with a practical yet provably secure construction requires some care. Below we give the details of our PRF implementation, which we call **TapeGen**.

### 5.1 The **TapeGen** PRF

**LENGTH-FLEXIBLE PRFS.** In practice, it is desirable that **TapeGen** be both variable input-length (VIL)- and variable output-length (VOL)-PRF,<sup>2</sup> a primitive we call a *length-flexible* (LF)-PRF. (In particular, the number of coins used by **HG** can be beyond one block of an underlying blockcipher in length, ruling out the use of most practical pseudorandom VIL-MACs.) That is, LF-PRF **TapeGen** with key-space  $Keys$  takes as input a key  $K \in Keys$ , an output length  $1^\ell$ , and  $x \in \{0, 1\}^*$  to return  $y \in \{0, 1\}^\ell$ . Define the following oracle  $R$  taking inputs  $1^\ell$  and  $x \in \{0, 1\}^*$  to return  $y \in \{0, 1\}^\ell$ , which maintains as state an array  $D$ :

**Oracle**  $R(1^\ell, x)$   
 If  $|D[x]| < \ell$  then  
      $r \xleftarrow{\$} \{0, 1\}^{\ell - |D[x]|}$   
      $D[x] \leftarrow D[x] \| r$   
 Return  $D[x]_1 \dots D[x]_\ell$

Above and in what follows,  $m_i$  denotes the  $i$ -th bit of a string  $m$ , and we require everywhere that  $\ell < \ell_{\max}$  for an associated maximum output length  $\ell_{\max}$ . For an adversary  $A$ , define its *length-flexible pseudorandom function (LF-PRF) advantage* against **TapeGen** as

$$\mathbf{Adv}_{\mathbf{TapeGen}}^{\text{prf}}(A) = \Pr \left[ A^{\mathbf{TapeGen}(K, \cdot, \cdot)} = 1 \right] - \Pr \left[ A^{R(\cdot, \cdot)} = 1 \right],$$

where the left probability is over the random choice of  $K \in Keys$ . Most practical VIL-MACs (message authentication codes) are PRFs and are therefore VIL-PRFs, but the VOL-PRF requirement does not seem to have been addressed previously. To achieve it we suggest using a VOL-PRG (pseudorandom generator) as well. Let us define the latter.

**VARIABLE-OUTPUT-LENGTH PRGS.** Let  $G$  be an algorithm that on input a seed  $s \in \{0, 1\}^k$  and an output length  $1^\ell$  returns  $y \in \{0, 1\}^\ell$ . Let  $\mathcal{O}_G$  be the oracle that on input  $1^\ell$  chooses a random seed  $s \in \{0, 1\}^k$  and returns  $G(s, \ell)$ , and let  $S$  be the oracle that on input  $1^\ell$  returns a random string  $r \in \{0, 1\}^\ell$ . For an adversary  $A$ , define its *variable-output-length pseudorandom function (VOL-PRG) advantage* against  $G$  as

$$\mathbf{Adv}_G^{\text{vol-prg}}(A) = \Pr \left[ A^{\mathcal{O}_G(\cdot)} = 1 \right] - \Pr \left[ A^{S(\cdot)} = 1 \right].$$

---

<sup>2</sup>That is, a VIL-PRF takes inputs of varying lengths. A VOL-PRF produces outputs of varying lengths specified by an additional input parameter.

As mentioned above, we require above that  $\ell < \ell_{\max}$  for an associated maximum output length  $\ell_{\max}$ . Call  $G$  *consistent* if  $\Pr [G(s, \ell') = G(s, \ell)_1 \dots G(s, \ell)_{\ell'}] = 1$  for all  $\ell' < \ell$ , with the probability over the choice of a random seed  $s \in \{0, 1\}^k$ . Most PRGs are consistent due to their “iterated” structure.

**OUR LF-PRF CONSTRUCTION.** We propose a general construction of an LF-PRF that composes a VIL-PRF with a consistent VOL-PRG by using the output of the former as the seed for the latter. Formally, let  $F$  be a VIL-PRF and  $G$  be a consistent VOL-PRG, and define the associated pseudorandom tape generation function **TapeGen** which on inputs  $K, 1^\ell, x$  returns  $G(1^\ell, F(K, x))$ .

The following says that **TapeGen** is indeed an LF-PRF if  $F$  is a VIL-PRF and  $G$  is a VOL-PRG.

**Proposition 5.1.** *Let  $A$  be an adversary against **TapeGen** that makes at most  $q$  queries to its oracle of total input length  $\ell_{\text{in}}$  and total output length  $\ell_{\text{out}}$ . Then there exists an adversary  $B_1$  against  $F$  and an adversary  $B_2$  against  $G$  such that*

$$\mathbf{Adv}_{\mathbf{TapeGen}}^{\text{prf}}(A) \leq \mathbf{Adv}_F^{\text{prf}}(B_1) + \mathbf{Adv}_G^{\text{vol-prg}}(B_2) .$$

*Adversaries  $B_1, B_2$  make at most  $q$  queries of total input length  $\ell_{\text{in}}$  and total output length  $\ell_{\text{out}}$  to their respective oracles and run in the time of  $A$ .*

*Proof.* We use a standard hybrid argument, changing the experiment where  $A$  has oracle  $\mathbf{TapeGen}(K, \cdot, \cdot)$  into one with oracle  $\mathcal{O}_R(\cdot, \cdot)$  in two steps. First change the former oracle to on input  $\ell, x$  output not  $G(\ell, F(K, x))$  but  $G(\ell, s)$  for a independent random  $s \in \{0, 1\}^k$ . The change in  $A$ 's advantage is bounded by  $\mathbf{Adv}_F^{\text{prf}}(B_1)$ , where  $B_1$  is the PRF adversary against  $F$  that runs  $A$ , responding to a query  $\ell, x$  by querying its own oracle with  $x$  to receive response  $y$ , and then returning  $G(\ell, y)$  to  $A$ . Next change  $A$ 's oracle to on input  $\ell, x$  return  $\mathcal{O}_R(\ell, x)$ . This time the change in  $A$ 's advantage is bounded by  $\mathbf{Adv}_G^{\text{vol-prg}}(B_2)$ , where  $B_2$  is the VOL-PRG adversary against  $G$  that runs  $A$ , responding to a query  $\ell, x$  with the response it receives to query  $\ell$  to its own oracle, and the proposition follows.  $\square$

Concretely, we suggest the following blockcipher-based consistent VOL-PRG for  $G$ . Let  $E: \{0, 1\}^k \times \{0, 1\}^n \rightarrow \{0, 1\}^n$  be a blockcipher. Define the associated VOL-PRG  $G[E]$  with seed-length  $k$  and maximum output length  $n \cdot 2^n$ , where  $G[E]$  on input  $s \in \{0, 1\}^k$  and  $1^\ell$  outputs the first  $\ell$  bits of  $E(s, \langle 1 \rangle) \| E(s, \langle 2 \rangle) \| \dots$  (Here  $\langle i \rangle$  denotes the  $n$ -bit binary encoding of  $i \in \mathbb{N}$ .) The following says that  $G[E]$  is a consistent VOL-PRG if  $E$  is a PRF.

**Proposition 5.2.** *Let  $E: \{0, 1\}^k \times \{0, 1\}^n \rightarrow \{0, 1\}^n$  be a blockcipher, and let  $A$  be an adversary against  $G[E]$  making at most  $q$  oracle queries whose responses total at most  $p \cdot n$  bits. Then there is an adversary  $B$  against  $E$  such that*

$$\mathbf{Adv}_{G[E]}^{\text{vol-prg}}(A) \leq q \cdot \mathbf{Adv}_E^{\text{prf}}(B) .$$

*Adversary  $B$  makes at most  $p$  queries to its oracle and runs in the time of  $A$ . Furthermore,  $G[E]$  is consistent.*

*Proof.* Consider the following adversary.

**Adversary**  $B^{\mathcal{O}(\cdot, \cdot)}$

$i \xleftarrow{\$} [q]$

$\text{ctr} \leftarrow 0$

Define  $\mathcal{P}$  as the oracle taking query  $1^\ell$  and running

$\text{ctr} \leftarrow \text{ctr} + 1$

If  $\text{ctr} < i$ :  $s \xleftarrow{\$} \{0, 1\}^k$  ; Return first  $\ell$  bits of  $E(s, \langle 1 \rangle) \| E(s, \langle 2 \rangle) \| \dots$

If  $\text{ctr} = i$ :  $s \xleftarrow{\$} \{0, 1\}^k$  ; Return first  $\ell$  bits of  $\mathcal{O}(s, \langle 1 \rangle) \| \mathcal{O}(s, \langle 2 \rangle) \| \dots$

If  $\text{ctr} > i$ :  $r \xleftarrow{\$} \{0, 1\}^\ell$  ; Return  $r$

$b \xleftarrow{\$} A^{\mathcal{P}(\cdot)}$

Return  $b$

In the PRF experiment,  $B$ 's oracle  $\mathcal{O}$  can be either the blockcipher  $E$  or a random function  $R : \{0, 1\}^k \times \{0, 1\}^n \rightarrow \{0, 1\}^n$ . Note that  $B$  with oracle  $E$  and  $i = 0$  emulates  $A$  with oracle  $G[E]$ ; while  $B$  with oracle  $R$  and  $i = q$  emulates  $A$  with oracle  $S$ , where  $S$  is the oracle that on input  $1^\ell$  returns a random string in  $\{0, 1\}^\ell$ . Hence,

$$\begin{aligned} \Pr \left[ A^{\mathcal{O}_{G[E]}(\cdot)} = 1 \right] &= \Pr \left[ B^{E(\cdot, \cdot)} = 1 \mid i = 1 \right], \\ \Pr \left[ A^S(\cdot) = 1 \right] &= \Pr \left[ B^{R(\cdot, \cdot)} = 1 \mid i = k \right]. \end{aligned} \quad (4)$$

Also, notice that for all  $j \in \{2, \dots, q-1\}$ ,  $B$  with oracle  $E$  and  $i = j$  has identical behavior to  $B$  with oracle  $R$  and  $i = j-1$ . Thus,

$$\Pr \left[ B^{E(\cdot, \cdot)} = 1 \mid i = j \right] = \Pr \left[ B^{R(\cdot, \cdot)} = 1 \mid i = j-1 \right] \quad \text{for all } j \in \{2, \dots, q\}. \quad (5)$$

Then,

$$\begin{aligned} &\mathbf{Adv}_{G[E]}^{\text{vol-prg}}(A) \\ &= \Pr \left[ A^{\mathcal{O}_{G[E]}(\cdot)} = 1 \right] - \Pr \left[ A^S(\cdot) = 1 \right] \\ &= \Pr \left[ B^{E(\cdot, \cdot)} = 1 \mid i = 1 \right] - \Pr \left[ B^{R(\cdot, \cdot)} = 1 \mid i = k \right] \quad [\text{by (4)}] \\ &= \sum_{j=1}^q \Pr \left[ B^{E(\cdot, \cdot)} = 1 \mid i = j \right] - \Pr \left[ B^{R(\cdot, \cdot)} = 1 \mid i = j \right] \quad [\text{by (5)}] \\ &= q \sum_{j=1}^q \Pr \left[ B^{E(\cdot, \cdot)} = 1 \mid i = j \right] \Pr[i = j] - \Pr \left[ B^{R(\cdot, \cdot)} = 1 \mid i = j \right] \Pr[i = j] \\ &= q \sum_{j=1}^q \Pr \left[ B^{E(\cdot, \cdot)} = 1 \cap i = j \right] - \Pr \left[ B^{R(\cdot, \cdot)} = 1 \cap i = j \right] \\ &\leq q \left( \Pr \left[ B^{E(\cdot, \cdot)} = 1 \right] - \Pr \left[ B^{R(\cdot, \cdot)} = 1 \right] \right) \\ &= q \mathbf{Adv}_E^{\text{prf}}(B) \end{aligned}$$

The efficiency claims should be clear from the definition of  $B$ . It is also obvious that  $G[E]$  is consistent: for  $\ell' < \ell$ , note that the first  $\ell'$  bits of  $G[E](s, 1^{\ell'})$  and  $G[E](s, 1^\ell)$  are the same as they are just the first  $\ell'$  bits of  $E(s, \langle 1 \rangle) \| E(s, \langle 2 \rangle) \| \dots$ .  $\square$

Now, to instantiate the VIL-PRF  $F$  in `TapeGen`, we suggest OMAC (a.k.a. CMAC) [22], which is also blockcipher-based and introduces no additional assumption. Then the secret key for `TapeGen` consists only of that for OMAC, which in turn consists of just one key for the underlying blockcipher (e.g. AES).

## 5.2 Our OPE Scheme and its Analysis

THE SCHEME. Let `TapeGen` be as above, with key-space  $Keys$ . Our associated order-preserving encryption scheme  $OPE^{HG}[\text{TapeGen}] = (\mathcal{K}^{HG}, \mathcal{E}nc^{HG}, \mathcal{D}ec^{HG})$  is defined as follows. The plaintext and ciphertext spaces are sets of consecutive integers  $\mathcal{D}, \mathcal{R}$ , respectively. Algorithm  $\mathcal{K}^{HG}$  returns a random  $K \in Keys$ . Algorithms  $\mathcal{E}nc^{HG}, \mathcal{D}ec^{HG}$  are the same as **LazySample**, **LazySampleInv**, respectively, except that HG is implemented by the algorithm of [23] and `GetCoins` by `TapeGen` (so there is no need to store the elements of  $F$  and  $I$ ). See Figure 2 for the formal descriptions of  $\mathcal{E}nc^{HG}$  and  $\mathcal{D}ec^{HG}$ , where as before  $\ell_1 = \ell(M, N, y - r)$  is the number of coins needed by HG on inputs  $M, N, y - r$ , and  $\ell_{\mathcal{R}}$  is the number of coins needed to select an element of  $\mathcal{R}$  uniformly at random. (The length parameters to `TapeGen` are just for convenience; one can always generate more output bits on-the-fly by invoking `TapeGen` again on a longer such parameter. In fact, our implementation of `TapeGen` can simply pick up where it left off instead of starting over.)

$\mathcal{E}nc_K^{HG}(\mathcal{D}, \mathcal{R}, m)$ 01 $M \leftarrow  \mathcal{D}  ; N \leftarrow  \mathcal{R} $ 02 $d \leftarrow \min(\mathcal{D}) - 1 ; r \leftarrow \min(\mathcal{R}) - 1$ 03 $y \leftarrow r + \lceil N/2 \rceil$ 04 If $ \mathcal{D}  = 1$ then 05 $cc \xleftarrow{\$} \text{TapeGen}(K, 1^{\ell_{\mathcal{R}}}, (\mathcal{D}, \mathcal{R}, 1  m))$ 06 $c \xleftarrow{cc} \mathcal{R}$ 07     Return $c$  08 $cc \xleftarrow{\$} \text{TapeGen}(K, 1^{\ell_1}, (\mathcal{D}, \mathcal{R}, 0  y))$ 09 $x \xleftarrow{\$} d + \text{HG}(M, N, y - r; cc)$ 10 If $m \leq x$ then 11 $\mathcal{D} \leftarrow \{d + 1, \dots, x\}$ 12 $\mathcal{R} \leftarrow \{r + 1, \dots, y\}$ 13 Else 14 $\mathcal{D} \leftarrow \{x + 1, \dots, d + M\}$ 15 $\mathcal{R} \leftarrow \{y + 1, \dots, r + N\}$ 16 Return $\mathcal{E}nc_K^{HG}(\mathcal{D}, \mathcal{R}, m)$	$\mathcal{D}ec_K^{HG}(\mathcal{D}, \mathcal{R}, c)$ 17 $M \leftarrow  \mathcal{D}  ; N \leftarrow  \mathcal{R} $ 18 $d \leftarrow \min(\mathcal{D}) - 1 ; r \leftarrow \min(\mathcal{R}) - 1$ 19 $y \leftarrow r + \lceil N/2 \rceil$ 20 If $ \mathcal{D}  = 1$ then $m \leftarrow \min(\mathcal{D})$ 21 $cc \xleftarrow{\$} \text{TapeGen}(K, 1^{\ell_{\mathcal{R}}}, (\mathcal{D}, \mathcal{R}, 1  m))$ 22 $w \xleftarrow{cc} \mathcal{R}$ 23     If $w = c$ then return $m$ 24     Else return $\perp$ 25 $cc \xleftarrow{\$} \text{TapeGen}(K, 1^{\ell_1}, (\mathcal{D}, \mathcal{R}, 0  y))$ 26 $x \xleftarrow{\$} d + \text{HG}(M, N, y - r; cc)$ 27 If $c \leq y$ then 28 $\mathcal{D} \leftarrow \{d + 1, \dots, x\}$ 29 $\mathcal{R} \leftarrow \{r + 1, \dots, y\}$ 30 Else 31 $\mathcal{D} \leftarrow \{x + 1, \dots, d + M\}$ 32 $\mathcal{R} \leftarrow \{y + 1, \dots, r + N\}$ 33 Return $\mathcal{D}ec_K^{HG}(\mathcal{D}, \mathcal{R}, c)$
---	--

Figure 2: Encryption  $\mathcal{E}nc^{HG}$  and decryption  $\mathcal{D}ec^{HG}$  algorithms for our hypergeometric distribution-based OPE scheme,  $OPE^{HG}[\text{TapeGen}]$ .

SECURITY. The following theorem characterizes security of our OPE scheme, saying that it is POPF-CCA secure if `TapeGen` is a LF-PRF. Applying Proposition 5.2, this is reduced to pseudo-randomness of an underlying blockcipher.

**Theorem 5.3.** Let  $\text{OPE}^{\text{HG}}[\text{TapeGen}]$  be the above OPE scheme with plaintext space size  $M$ , ciphertext space size  $N$ . For adversary  $A$  against  $\text{OPE}^{\text{HG}}[\text{TapeGen}]$  making at most  $q$  queries to its oracles combined, there is an adversary  $B$  against  $\text{TapeGen}$  such that

$$\text{Adv}_{\text{OPE}^{\text{HG}}[\text{TapeGen}]}^{\text{popf-cca}}(A) \leq \text{Adv}_{\text{TapeGen}}^{\text{prf}}(B) + \lambda.$$

Adversary  $B$  makes at most  $q_1 = q \cdot (\log N + 1)$  queries of size at most  $5 \log N + 1$  to its oracle, whose responses total  $q_1 \cdot \lambda'$  bits on average, and its running-time is that of  $A$ . Above,  $\lambda, \lambda'$  are constants depending only on  $\text{HG}$  and the precision of the underlying floating-point computations (not on  $M, N$ ).

*Proof.* Define adversary  $B$  as follows. Given an oracle for either  $\text{TapeGen}$  or a random function with corresponding inputs and outputs lengths,  $B$  runs  $A$  and replies to its oracle queries by simulating  $\mathcal{Enc}^{\text{HG}}$  and  $\mathcal{Dec}^{\text{HG}}$  algorithms. Note that only the procedure  $\text{TapeGen}$  used by these algorithms uses the secret key.  $B$  simulates it using its own oracle. We have

$$\begin{aligned} & \text{Adv}_{\text{OPE}^{\text{HG}}[\text{TapeGen}]}^{\text{popf-cca}}(A) \\ &= \Pr \left[ A^{\mathcal{Enc}^{\text{HG}}(K, \cdot), \mathcal{Dec}^{\text{HG}}(K, \cdot)} = 1 \right] - \Pr \left[ A^{g(\cdot), g^{-1}(\cdot)} = 1 \right] \\ &= \Pr \left[ A^{\mathcal{Enc}^{\text{HG}}(K, \cdot), \mathcal{Dec}^{\text{HG}}(K, \cdot)} = 1 \right] - \Pr \left[ A^{\text{LazySample}(\mathcal{D}, \mathcal{R}, \cdot), \text{LazySampleInv}(\mathcal{D}, \mathcal{R}, \cdot)} = 1 \right] \\ &\leq \text{Adv}_{\text{TapeGen}}^{\text{prf}}(B) + \lambda. \end{aligned}$$

The first equation is by definition. The second equation is due to Theorem 4.2. The last inequality is justified as follows. By construction our  $\mathcal{Enc}^{\text{HG}}$  and  $\mathcal{Dec}^{\text{HG}}$  algorithms differ from **LazySample** and **LazySampleInv** respectively only in the use of random tape, which is truly random in one case and pseudorandom in another. Thus any difference in the probabilities in the second line will equal the difference  $B$ 's output distribution which is  $\text{Adv}_{\text{TapeGen}}^{\text{prf}}(B)$ . Above  $\lambda$  represents an “error term” due to the fact that the “exact” hypergeometric sampling algorithm of [23] technically requires infinite floating-point precision, which is not possible in the real world. One way to bound  $\lambda$  would be to bound the probability that an adversary can distinguish the used  $\text{HG}$  sampling algorithm from the ideal (infinite precision) one.  $B$ 's running time and resources are justified by observing the algorithms and their efficiency analysis.  $\square$

**EFFICIENCY.** The efficiency of our scheme follows from our previous analyses. Using the suggested implementation of  $\text{TapeGen}$  in Subsection 5.1, encryption and decryption require the time for at most  $\log N + 1$  invocations of  $\text{HG}$  on inputs of size at most  $\log N$  plus at most  $(5 \log M + 12) \cdot (5 \log N + \lambda' + 1)/128$  invocations of AES on average for  $\lambda'$  in the theorem.

### 5.3 On Choosing $N$

Practitioners interested in implementing our scheme might naturally wonder how large we recommend making the ciphertext space size  $N$ . In fact, different choices of  $N$  have no bearing on our scheme's POPF-CCA security. Rather, different choices of  $N$  will affect how the ideal object, a random OPF, behaves. Thus, in order to say something meaningful about the choice of  $N$ , we first need a security definition and analysis for the ideal object, which is a topic of ongoing research.

## 6 On Using the Negative Hypergeometric Distribution

In the balls-and-bins model described in Section 4.1 with  $M$  black and  $N - M$  white balls in the bin, consider the random variable  $Y$  describing the total number of balls in our sample after we pick the  $x$ -th black ball. This random variable follows the *negative* hypergeometric (NHG) distribution. Formally,

$$P_{\text{NHG}}(y; N, M, x) = \frac{\binom{y-1}{x-1} \cdot \binom{N-y}{M-x}}{\binom{N}{M}}.$$

As we discussed in the introduction, use of the NHG distribution instead of the HG permits slightly simpler and more efficient lazy sampling algorithms and corresponding OPE scheme. The problem is that they require an efficient NHG sampling algorithm, and the existence of such an algorithm is apparently open. What is known is that the NHG distribution can be approximated by the negative binomial distribution [27], the latter can be sampled efficiently [17, 15], and the approximation improves as  $M$  and  $N$  grow. However, quantifying the quality of the approximation for fixed parameters seems difficult. If future work either develops an efficient exact sampling algorithm for the NHG distribution or shows that the approximation by the negative binomial distribution is sufficiently close, then our NHG-based OPE scheme could be a good alternative to the HG-based one. Here are the details.

### 6.1 Construction of the NHG-based OPE Scheme

Assume there exists an efficient algorithm  $\text{NHG}$  that efficiently samples according to the NHG distribution, possibly using an approximation to a related distribution as we discussed.  $\text{NHG}$  takes inputs  $M, N$ , and  $x \in \{0, 1, \dots, M\}$  and returns  $y \in \{0, 1, \dots, N\}$  such that for each  $y^* \in \{0, 1, \dots, N\}$  we have  $y = y^*$  with probability  $P_{\text{NHG}}(y^*; N, M, x)$  over the coins of  $\text{NHG}$ . Let  $\ell_1 = \ell(M, N, y - r)$  denote the number of coins needed by  $\text{NHG}$  on inputs  $M, N, y - r$ .

Define  $\text{OPE}^{\text{NHG}}[\text{TapeGen}] = (\mathcal{K}, \text{Enc}^{\text{NHG}}, \text{Dec}^{\text{NHG}})$ , our NHG-based order-preserving encryption scheme, as follows. Let  $\text{TapeGen}$  be the PRF described in Section 5, with key-space  $\text{Keys}$ . The plaintext and ciphertext spaces are sets of consecutive integers  $\mathcal{D}, \mathcal{R}$ , respectively. Algorithm  $\mathcal{K}$  returns a random  $K \in \text{Keys}$ . Algorithms  $\text{Enc}^{\text{NHG}}, \text{Dec}^{\text{NHG}}$  are described in Figure 3.

### 6.2 Correctness

We prove correctness of the NHG scheme in the same manner as the HG scheme. First, see in Figure 4 the revised versions  $\text{LazySample}^*$ ,  $\text{LazySampleInv}^*$  of the stateful algorithms from before. The algorithms re-use the subroutine  $\text{GetCoins}$ , which takes inputs  $1^\ell, \mathcal{D}, \mathcal{R}$ , and  $b||z$ , where  $b \in \{0, 1\}$  and  $z \in \mathcal{R}$  if  $b = 0$  and  $z \in \mathcal{D}$  otherwise, to return  $cc \in \{0, 1\}^\ell$ . Also, recall that the array  $I$ , initially empty, is global and shared between the algorithms.

With these revised versions of  $\text{LazySample}^*$ ,  $\text{LazySampleInv}^*$ , we supply a revised version of Theorem 4.2 for the NHG case.

**Theorem 6.1.** *Suppose  $\text{GetCoins}$  returns truly random coins on each new input. Then for any (even computationally unbounded) algorithm  $A$  we have*

$$\Pr \left[ A^{g(\cdot), g^{-1}(\cdot)} = 1 \right] = \Pr \left[ A^{\text{LazySample}^*(\mathcal{D}, \mathcal{R}, \cdot), \text{LazySampleInv}^*(\mathcal{D}, \mathcal{R}, \cdot)} = 1 \right],$$

$\mathcal{E}nc_K^{\text{NHG}}(\mathcal{D}, \mathcal{R}, m)$ 16 If $ \mathcal{D}  = 0$ then return $\perp$ 17 $M \leftarrow  \mathcal{D} $ ; $N \leftarrow  \mathcal{R} $ 18 $d \leftarrow \min(\mathcal{D}) - 1$ 19 $r \leftarrow \min(\mathcal{R}) - 1$ 20 $x \leftarrow d + \lceil M/2 \rceil$ 21 $cc \stackrel{\$}{\leftarrow} \text{TapeGen}(K, 1^{\ell_1}, (\mathcal{D}, \mathcal{R}, x))$ 22 $y \leftarrow r + \text{NHG}(N, M, x - d; cc)$ 23 If $m = x$ then 24     Return $y$ 25 If $m < x$ then 26 $\mathcal{D} \leftarrow \{d + 1, \dots, x - 1\}$ 27 $\mathcal{R} \leftarrow \{r + 1, \dots, y - 1\}$ 28 Else 29 $\mathcal{D} \leftarrow \{x + 1, \dots, d + M\}$ 30 $\mathcal{R} \leftarrow \{y + 1, \dots, r + N\}$ 31 Return $\mathcal{E}nc_K^{\text{NHG}}(\mathcal{D}, \mathcal{R}, m)$	$\mathcal{D}ec_K^{\text{NHG}}(\mathcal{D}, \mathcal{R}, c)$ 16 If $ \mathcal{D}  = 0$ then return $\perp$ 17 $M \leftarrow  \mathcal{D} $ ; $N \leftarrow  \mathcal{R} $ 18 $d \leftarrow \min(\mathcal{D}) - 1$ 19 $r \leftarrow \min(\mathcal{R}) - 1$ 20 $x \leftarrow d + \lceil M/2 \rceil$ 21 $cc \stackrel{\$}{\leftarrow} \text{TapeGen}(K, 1^{\ell_1}, (\mathcal{D}, \mathcal{R}, x))$ 22 $y \leftarrow r + \text{NHG}(N, M, x - d; cc)$ 23 If $c = y$ then 24     Return $x$ 25 If $c < y$ then 26 $\mathcal{D} \leftarrow \{d + 1, \dots, x - 1\}$ 27 $\mathcal{R} \leftarrow \{r + 1, \dots, y - 1\}$ 28 Else 29 $\mathcal{D} \leftarrow \{x + 1, \dots, d + M\}$ 30 $\mathcal{R} \leftarrow \{y + 1, \dots, r + N\}$ 31 Return $\mathcal{D}ec_K^{\text{NHG}}(\mathcal{D}, \mathcal{R}, c)$
---	---

Figure 3: Encryption  $\mathcal{E}nc^{\text{NHG}}$  and decryption  $\mathcal{D}ec^{\text{NHG}}$  algorithms for our negative hypergeometric distribution-based OPE scheme,  $\text{OPE}^{\text{NHG}}[\text{TapeGen}]$ .

where  $g, g^{-1}$  denote an order-preserving function picked at random from  $\text{OPF}_{\mathcal{D}, \mathcal{R}}$  and its inverse, respectively.

*Proof.* Since we consider unbounded adversaries, we can ignore the inverse oracle in our analysis, since such an adversary can always query all points in the domain to learn all points in the image. Let  $M = |\mathcal{D}|$ ,  $N = |\mathcal{R}|$ ,  $d = \min(\mathcal{D}) - 1$ , and  $r = \min(\mathcal{R}) - 1$ . We will say that two functions  $g, h : \mathcal{D} \rightarrow \mathcal{R}$  are *equivalent* if  $g(m) = h(m)$  for all  $m \in \mathcal{D}$ . (Note that if  $\mathcal{D} = \emptyset$ , any two functions  $g, h : \mathcal{D} \rightarrow \mathcal{R}$  are vacuously equivalent.) Let  $f$  be any function in  $\text{OPF}_{\mathcal{D}, \mathcal{R}}$ . To prove the theorem, it is enough to show that the function defined by **LazySample** $^*(\mathcal{D}, \mathcal{R}, \cdot)$  is equivalent to  $f$  with probability  $1/|\text{OPF}_{\mathcal{D}, \mathcal{R}}|$ . We prove this using strong induction on  $M$  and  $N$ .

Consider the base case where  $M = 1$ , i.e.,  $\mathcal{D} = \{m\}$  for some  $m$ , and  $N \geq M$ . When it is first called, **LazySample** $^*(\mathcal{D}, \mathcal{R}, m)$  will determine random coins  $cc$ , then enter the result of  $\text{NHG}(M, N, m - d; cc)$  into  $I[\mathcal{D}, \mathcal{R}, m]$ , whereupon this any future calls of **LazySample** $^*(\mathcal{D}, \mathcal{R}, m)$  will always output  $F[\mathcal{D}, \mathcal{R}, m] = c$ . Note that by definition,  $\text{NHG}(M, N, m - d; cc)$  returns  $f(m)$  with probability

$$P_{\text{NHG}}(f(m) - r; N, 1, 1) = \frac{\binom{f(m) - r - 1}{0} \cdot \binom{N - (f(m) - r)}{0}}{\binom{N}{1}} = \frac{1}{N} = \frac{1}{|\mathcal{R}|}.$$

Thus, the output of **LazySample** $^*(\mathcal{D}, \mathcal{R}, m)$  will always be  $f(m)$  with probability  $1/|\mathcal{R}|$ , implying that **LazySample** $^*(\mathcal{D}, \mathcal{R}, m)$  is equivalent to  $f(m)$  with probability  $1/|\mathcal{R}| = 1/|\text{OPF}_{\mathcal{D}, \mathcal{R}}|$ .

Now suppose  $M > 1$ , and  $N \geq M$ . As an induction hypothesis assume that for all domains  $\mathcal{D}'$  of size  $M'$  and ranges  $\mathcal{R}'$  of size  $N' \geq M'$ , where either  $M' < M$  or ( $M' = M$  and  $N' < N$ ), and for any function  $f'$  in  $\text{OPF}_{\mathcal{D}', \mathcal{R}'}$ , **LazySample** $^*(\mathcal{D}', \mathcal{R}', \cdot)$  is equivalent to  $f'$  with probability  $1/|\text{OPF}_{\mathcal{D}', \mathcal{R}'}|$ .

<b>LazySample*</b> ( $\mathcal{D}, \mathcal{R}, m$ )	<b>LazySampleInv*</b> ( $\mathcal{D}, \mathcal{R}, c$ )
01 $M \leftarrow  \mathcal{D} ; N \leftarrow  \mathcal{R} $	17 If $ \mathcal{D}  = 0$ then return $\perp$
02 $d \leftarrow \min(\mathcal{D}) - 1; r \leftarrow \min(\mathcal{R}) - 1$	18 $M \leftarrow  \mathcal{D} ; N \leftarrow  \mathcal{R} $
03 $x \leftarrow d + \lceil M/2 \rceil$	19 $d \leftarrow \min(\mathcal{D}) - 1; r \leftarrow \min(\mathcal{R}) - 1$
04 If $I[\mathcal{D}, \mathcal{R}, x]$ is undefined then	20 $x \leftarrow d + \lceil M/2 \rceil$
05 $cc \xleftarrow{\$} \text{GetCoins}(1^{\ell_1}, \mathcal{D}, \mathcal{R}, 1  x)$	21 If $I[\mathcal{D}, \mathcal{R}, x]$ is undefined then
06 $I[\mathcal{D}, \mathcal{R}, x] \xleftarrow{\$}$ $\text{NHG}(M, N, x - d; cc)$	22 $cc \xleftarrow{\$} \text{GetCoins}(1^{\ell_1}, \mathcal{D}, \mathcal{R}, 1  x)$
07 $y \leftarrow r + I[\mathcal{D}, \mathcal{R}, x]$	23 $I[\mathcal{D}, \mathcal{R}, x] \xleftarrow{\$}$ $\text{NHG}(M, N, x - d; cc)$
08 If $m = x$ then	24 $y \leftarrow r + I[\mathcal{D}, \mathcal{R}, x]$
09 Return $y$	25 If $c = y$ then
10 If $m < x$ then	26 Return $x$
11 $\mathcal{D} \leftarrow \{d + 1, \dots, x - 1\}$	27 If $c < y$ then
12 $\mathcal{R} \leftarrow \{r + 1, \dots, y - 1\}$	28 $\mathcal{D} \leftarrow \{d + 1, \dots, x - 1\}$
13 Else	29 $\mathcal{R} \leftarrow \{r + 1, \dots, y - 1\}$
14 $\mathcal{D} \leftarrow \{x + 1, \dots, d + M\}$	30 Else
15 $\mathcal{R} \leftarrow \{y + 1, \dots, r + N\}$	31 $\mathcal{D} \leftarrow \{x + 1, \dots, d + M\}$
16 Return <b>LazySample*</b> ( $\mathcal{D}, \mathcal{R}, m$ )	32 $\mathcal{R} \leftarrow \{y + 1, \dots, r + N\}$
	33 Return <b>LazySampleInv*</b> ( $\mathcal{D}, \mathcal{R}, c$ )

Figure 4: Algorithms **LazySample\*** and **LazySampleInv\*** for lazy-sampling a pseudorandom order-preserving function and its inverse by sampling the negative hypergeometric distribution.

The first time it is called, **LazySample\***( $\mathcal{D}, \mathcal{R}, \cdot$ ) first computes  $I[\mathcal{D}, \mathcal{R}, x]$  from  $\text{NHG}(M, N, x - d; cc)$ , where  $x = d + \lceil M/2 \rceil$ . Henceforth, on this and future calls of **LazySample\***( $\mathcal{D}, \mathcal{R}, \cdot$ ), the algorithm sets  $y \leftarrow r + I[\mathcal{D}, \mathcal{R}, x]$ , and follows one of three routes: if  $x = m$ , the algorithm terminates and returns  $y$ , if  $m < x$  it will return the output of **LazySample\***( $\mathcal{D}_1, \mathcal{R}_1, m$ ), and if  $m > x$  it will return the output of **LazySample\***( $\mathcal{D}_2, \mathcal{R}_2, m$ ), where  $\mathcal{D}_1 = \{1, \dots, x - 1\}$ ,  $\mathcal{R}_1 = \{1, \dots, y - 1\}$ ,  $\mathcal{D}_2 = \{x + 1, \dots, M\}$ ,  $\mathcal{R}_2 = \{y + 1, \dots, N\}$ . Let  $f_1$  be  $f$  restricted to the domain  $\mathcal{D}_1$ , and let  $f_2$  be  $f$  restricted to the domain  $\mathcal{D}_2$ . Note then that **LazySample\***( $\mathcal{D}, \mathcal{R}, \cdot$ ) is equivalent to  $f$  if and only if all three of the following events occur:

$E_1$ : The invocation of  $\text{NHG}(M, N, x - d; cc)$  returns the value  $f(x) - r$ .

$E_2$ : **LazySample\***( $\mathcal{D}_1, \mathcal{R}_1, \cdot$ ) is equivalent to  $f_1$ .

$E_3$ : **LazySample\***( $\mathcal{D}_2, \mathcal{R}_2, \cdot$ ) is equivalent to  $f_2$ .

By the law of conditional probability, and since  $E_2$  and  $E_3$  are independent,

$$\begin{aligned} \Pr[E_1 \cap E_2 \cap E_3] &= \Pr[E_1] \Pr[E_2 \cap E_3 \mid E_1] \\ &= \Pr[E_1] \Pr[E_2 \mid E_1] \Pr[E_3 \mid E_1]. \end{aligned}$$

$\Pr[E_1]$  is the negative hypergeometric probability that  $\text{NHG}(M, N, x - d)$  will return  $f(x) - r$ , which is

$$\Pr[E_1] = P_{\text{NHG}}(f(x) - r; N, M, \lceil M/2 \rceil) = \frac{\binom{f(x) - r - 1}{\lceil M/2 \rceil - 1} \binom{N - f(x) + r}{M - \lceil M/2 \rceil}}{\binom{N}{M}}.$$



Assume that  $E_1$  holds, and thus  $f_1$  is an element of  $\text{OPF}_{\mathcal{D}_1, \mathcal{R}_1}$  and  $f_2$  is an element of  $\text{OPF}_{\mathcal{D}_2, \mathcal{R}_2}$ . By definition,  $|\mathcal{R}_1|, |\mathcal{R}_2| < |\mathcal{R}|$ , and  $|\mathcal{D}_1|, |\mathcal{D}_2| \leq |\mathcal{D}|$ . So the induction hypothesis holds for each, and thus  $\mathbf{LazySample}^*(\mathcal{D}_1, \mathcal{R}_1, \cdot)$  is equivalent to  $f_1$  with probability  $1/|\text{OPF}_{\mathcal{D}_1, \mathcal{R}_1}| = 1/\binom{|\mathcal{R}_1|}{|\mathcal{D}_1|}$ , and  $\mathbf{LazySample}^*(\mathcal{D}_2, \mathcal{R}_2, \cdot)$  is equivalent to  $f_2$  with probability  $1/|\text{OPF}_{\mathcal{D}_2, \mathcal{R}_2}| = 1/\binom{|\mathcal{R}_2|}{|\mathcal{D}_2|}$ . Thus, we have that

$$\Pr[E_2 \mid E_1] = \frac{1}{\binom{f(x)-r-1}{\lceil M/2 \rceil - 1}} \quad \text{and} \quad \Pr[E_3 \mid E_1] = \frac{1}{\binom{N-f(x)+r}{M - \lceil M/2 \rceil}}.$$

Thus,

$$\Pr[E_1 \cap E_2 \cap E_3] = \frac{\binom{f(x)-r-1}{\lceil M/2 \rceil - 1} \binom{N-f(x)+r}{M - \lceil M/2 \rceil}}{\binom{N}{M}} \frac{1}{\binom{f(x)-r-1}{\lceil M/2 \rceil - 1}} \frac{1}{\binom{N-f(x)+r}{M - \lceil M/2 \rceil}} = \frac{1}{\binom{N}{M}}.$$

Therefore,  $\mathbf{LazySample}^*(\mathcal{D}, \mathcal{R}, \cdot)$  is equivalent to  $f$  with probability  $\frac{1}{\binom{N}{M}} = \frac{1}{|\text{OPF}_{\mathcal{D}, \mathcal{R}}|}$ . Since  $f$  was an arbitrary element of  $\text{OPF}_{\mathcal{D}, \mathcal{R}}$ , the result follows.  $\square$

Now, it is straightforward to prove the formal statement of correctness as before.

**Theorem 6.2.** *Let  $\text{OPE}^{\text{NHG}}[\text{TapeGen}]$  be the OPE scheme defined above with plaintext-space of size  $M$  and ciphertext space of size  $N$ . Then for any adversary  $A$  against  $\text{OPE}^{\text{NHG}}[\text{TapeGen}]$  making at most  $q$  queries to its oracles combined, there is an adversary  $B$  against  $\text{TapeGen}$  such that*

$$\mathbf{Adv}_{\text{OPE}^{\text{NHG}}[\text{TapeGen}]}^{\text{popf-cca}}(A) \leq \mathbf{Adv}_{\text{TapeGen}}^{\text{prf}}(B) + \lambda.$$

*Adversary  $B$  makes at most  $q_1 = q \cdot (\log N + 1)$  queries of size at most  $5 \log N + 1$  to its oracle, whose responses total  $q_1 \cdot \lambda'$  bits on average, and its running-time is that of  $A$ . Above,  $\lambda, \lambda'$  are constants depending only on NHG and the precision of the underlying floating-point computations (not on  $M, N$ ).*

*Proof.* The proof of this theorem is identical to that of Theorem 5.3, except that it uses Theorem 6.1 as a lemma rather than Theorem 4.2.  $\square$

### 6.3 Efficiency of the NHG Scheme

Efficiency-wise, it is not hard to see that to encrypt a single plaintext, each algorithm performs  $\log M + 1$  recursions in the worst-case (as opposed to  $\log N + 1$  for the HG-based algorithms), as the algorithm finds the desired plaintext via a binary search over the plaintext space, at each recursion calling NHG to determine the encryption of the midpoint (defined as the last plaintext in the first half of the current plaintext domain). The expected number of recursions is easily deduced as

$$\frac{1}{M} \cdot \left[ (\log M + 1) + \sum_{k=1}^{\log M} 2^{k-1} k \right].$$

A simple inductive proof shows that this value is between  $\log M - 1$  and  $\log M$ . This falls in line with what we expect from a binary-search strategy, where the expected number of iterations is typically only about 1 fewer than the worst-case number of iterations.

The algorithms of the corresponding OPE scheme can be obtained following the same idea of eliminating state by using a length-flexible PRF as described in Section 5.2. The security statement is the same as that of Theorem 5.3, where the last term now corresponds to the error probability of the NHG algorithm.

## Acknowledgements

We thank Anna Lysyanskaya, Silvio Micali, Leonid Reyzin, Ron Rivest, Phil Rogaway and the anonymous reviewers of Eurocrypt 2009 for helpful comments and references. Alexandra Boldyreva, Nathan Chenette and Adam O’Neill were supported in part by Alexandra’s NSF CAREER award 0545659 and NSF Cyber Trust award 0831184. Younho Lee was supported in part by the Korea Research Foundation Grant funded by the Korean Government (MOEHRD) (KRF:2007-357-D00243). Also, he is supported by Professor Mustaque Ahamad through the funding provided by IBM ISS and AT&T.

## References

- [1] R. Agrawal, J. Kiernan, R. Srikant, and Y. Xu. Order-preserving encryption for numeric data. In G. Weikum, A. C. König, and S. Deßloch, editors, *SIGMOD Conference*, pages 563–574. ACM, 2004. 3
- [2] G. Amanatidis, A. Boldyreva, and A. O’Neill. Provably-secure schemes for basic query support in outsourced databases. In S. Barker and G.-J. Ahn, editors, *DBSec*, volume 4602 of *Lecture Notes in Computer Science*, pages 14–30. Springer, 2007. 3
- [3] F. L. Bauer. *Decrypted Secrets: Methods and Maxims of Cryptology*. Springer, 2006. 3
- [4] M. Bellare. New proofs for NMAC and HMAC: Security without collision-resistance. In C. Dwork, editor, *CRYPTO*, volume 4117 of *Lecture Notes in Computer Science*, pages 602–619. Springer, 2006. 4
- [5] M. Bellare, A. Boldyreva, L. R. Knudsen, and C. Namprempe. Online ciphers and the Hash-CBC construction. In J. Kilian, editor, *CRYPTO*, volume 2139 of *Lecture Notes in Computer Science*, pages 292–309. Springer, 2001. 3, 4, 10, 11
- [6] M. Bellare, A. Boldyreva, and A. O’Neill. Deterministic and efficiently searchable encryption. In A. Menezes, editor, *CRYPTO*, volume 4622 of *Lecture Notes in Computer Science*, pages 535–552. Springer, 2007. 3, 4, 6
- [7] M. Bellare, M. Fischlin, A. O’Neill, and T. Ristenpart. Deterministic encryption: Definitional equivalences and constructions without random oracles. In D. Wagner, editor, *CRYPTO*, volume 5157 of *Lecture Notes in Computer Science*, pages 360–378. Springer, 2008. 3, 4
- [8] M. Bellare, T. Kohno, and C. Namprempe. Breaking and provably repairing the SSH authenticated encryption scheme: A case study of the Encode-then-Encrypt-and-MAC paradigm. *ACM Trans. Inf. Syst. Secur.*, 7(2):206–241, May 2004. 4, 8

- [9] M. Bellare and P. Rogaway. The security of triple encryption and a framework for code-based game-playing proofs. In *EUROCRYPT*, Lecture Notes in Computer Science, pages 409–426. Springer, 2006. 4, 11
- [10] A. Boldyreva, N. Chenette, Y. Lee, and A. O’Neill. Order-preserving symmetric encryption. In A. Joux, editor, *EUROCRYPT*, volume 5479 of *Lecture Notes in Computer Science*, pages 224–241. Springer, 2009. 9
- [11] A. Boldyreva, S. Fehr, and A. O’Neill. On notions of security for deterministic encryption, and efficient constructions without random oracles. In D. Wagner, editor, *CRYPTO*, volume 5157 of *Lecture Notes in Computer Science*, pages 335–359. Springer, 2008. 3, 4
- [12] D. Boneh and B. Waters. Conjunctive, subset, and range queries on encrypted data. In S. P. Vadhan, editor, *TCC*, volume 4392 of *Lecture Notes in Computer Science*, pages 535–554. Springer, 2007. 3
- [13] V. Chvátal. The Tail of the Hypergeometric Distribution. *Discrete Mathematics*, 25:285–287, 1979. 16
- [14] Z. Erkin, A. Piva, S. Katzenbeisser, R. L. Legendijk, J. Shokrollahi, G. Neven, and M. Barni. Protection and retrieval of encrypted multimedia content: When cryptography meets signal processing. *EURASIP J. Information Security*, 2007. 3
- [15] G. Fishman. *Discrete-Event Simulation: Modeling, Programming, and Analysis*. Springer Series in Operations Research. Springer, 2001. 5, 16, 22
- [16] E. A. Fox, Q. F. Chen, A. M. Daoud, and L. S. Heath. Order-preserving minimal perfect hash functions and information retrieval. *ACM Trans. Inf. Syst.*, 9(3):281–308, July 1991. 6
- [17] J. Gentle. *Random Number Generation and Monte Carlo Methods*. Statistics and Computing. Springer, 2003. 5, 22
- [18] O. Goldreich, S. Goldwasser, and S. Micali. How to construct random functions. *J. ACM*, 33(4):792–807, Aug. 1986. 10
- [19] O. Goldreich, S. Goldwasser, and A. Nussboim. On the implementation of huge random objects. *SIAM J. Comput.*, 39(7):2761–2822, 2010. 4
- [20] L. Granboulan and T. Pornin. Perfect block ciphers with small blocks. In A. Biryukov, editor, *FSE*, volume 4593 of *Lecture Notes in Computer Science*, pages 452–465. Springer, 2007. 6
- [21] P. Indyk, R. Motwani, P. Raghavan, and S. Vempala. Locality-preserving hashing in multidimensional spaces. In F. T. Leighton and P. W. Shor, editors, *STOC*, pages 618–625. ACM, 1997. 6
- [22] T. Iwata and K. Kurosawa. OMAC: One-key CBC MAC. In T. Johansson, editor, *FSE*, volume 2887 of *Lecture Notes in Computer Science*, pages 129–153. Springer, 2003. 20
- [23] V. Kachitvichyanukul and B. Schmeiser. Computer generation of hypergeometric random variates. *Statistical Computation and Simulation*, 22:127–145, 1985. 5, 6, 16, 20, 21

- [24] V. Kachitvichyanukul and B. W. Schmeiser. Algorithm 668: H2PEC: sampling from the hypergeometric distribution. *ACM Trans. Math. Softw.*, 14(4):397–398, 1988. 5, 6, 16
- [25] J. Li and E. Omiecinski. Efficiency and security trade-off in supporting range queries on encrypted databases. In S. Jajodia and D. Wijesekera, editors, *DBSec*, volume 3654 of *Lecture Notes in Computer Science*, pages 69–83. Springer, 2005. 3
- [26] N. Linial and O. Sasson. Non-expansive hashing. In G. L. Miller, editor, *STOC*, pages 509–518. ACM, 1996. 6
- [27] F. López-Blázquez and B. Salamanca-Miño. Exact and approximated relations between negative hypergeometric and negative binomial probabilities. *Communications in Statistics - Theory and Methods*, 30(5):957–967, 2001. 5, 22
- [28] P. Rogaway and T. Shrimpton. Deterministic authenticated-encryption: A provable-security treatment of the key-wrap problem. *IACR Cryptology ePrint Archive*, 2006. 4
- [29] A. C. C. Say and A. K. Nircan. Random generation of monotonic functions for Monte Carlo solution of qualitative differential equations. *Automatica*, 41(5):739–754, 2005. 11
- [30] E. Shi, J. Bethencourt, H. T.-H. Chan, D. X. Song, and A. Perrig. Multi-dimensional range query over encrypted data. In *IEEE Symposium on Security and Privacy*, pages 350–364. IEEE Computer Society, 2007. 3
- [31] A. J. Walker. An efficient method for generating discrete random variables with general distributions. *ACM Transactions on Mathematical Software*, 3(3):253–256, 1977. 16
- [32] D. Westhoff, J. Girão, and M. Acharya. Concealed data aggregation for reverse multicast traffic in sensor networks: Encryption, key distribution, and routing adaptation. *IEEE Trans. Mob. Comput.*, 5(10):1417–1431, 2006. 3
- [33] J. Xu, J. Fan, M. H. Ammar, and S. B. Moon. Prefix-preserving IP address anonymization: Measurement-based security evaluation and a new cryptography-based scheme. In *ICNP*, pages 280–289. IEEE Computer Society, 2002. 3