

# Skye: An Expanding PRF based Fast KDF and its Applications

Amit Singh Bhati  
COSIC, KU Leuven  
Leuven, Belgium  
amitsingh.bhati@esat.kuleuven.be

Antonín Dufka  
Masaryk University  
Brno, Czech Republic  
dufkan@mail.muni.cz

Elena Andreeva  
Technical University of Vienna  
Vienna, Austria  
elena.andreeva@tuwien.ac.at

Arnab Roy  
Alpen-Adria University  
Klagenfurt, Austria  
arnab.roy@aau.at

Bart Preneel  
COSIC, KU Leuven  
Leuven, Belgium  
bart.preneel@esat.kuleuven.be

## ABSTRACT

A Key Derivation Function (KDF) generates a uniform and highly random key-stream from weakly random key material. KDFs are broadly used in various security protocols such as digital signatures and key exchange protocols. HKDF, the most deployed KDF in practice, is based on the *extract-then-expand* paradigm. It is presently used, among others, in the Signal Protocol for end-to-end encrypted messaging.

HKDF is a generic KDF for general input sources and thus is not optimized for source-specific use cases such as key derivation from Diffie-Hellman (DH) sources (i.e. DH shared secrets as key material). Furthermore, the sequential HKDF design is unnecessarily slow on some general-purpose platforms that can benefit from parallelization.

In this work, we propose a novel, efficient and secure KDF called Skye. Skye follows the *extract-then-expand* paradigm and consists of two algorithms: efficient deterministic *randomness extractor* and *expansion* functions. Instantiating our extractor for dedicated source-specific (e.g. DH sources) inputs leads to a significant efficiency gain over HKDF while maintaining its security level. We provide concrete security analysis of Skye and both its underlying algorithms in the standard model.

We provide a software performance comparison of Skye with the AES-based expanding PRF ButterKnife and HKDF with SHA-256 (as used in practice). Our results show that in isolation Skye performs from 4x to 47x faster than HKDF, depending on the availability of AES or SHA instruction support. We further demonstrate that with such a performance gain, when Skye is integrated within the current Signal implementation, we can achieve significant overall improvements ranging from 38% to 64% relative speedup in unidirectional messaging. Even in bidirectional messaging, that includes DH computation with dominating computational cost, Skye still contributes to 12-36% relative speedup when just 10 messages are sent and received at once.

## 1 INTRODUCTION

A Cryptographic Key Derivation Function (KDF) outputs a uniform and “highly” random arbitrarily long key-stream when provided with a non-uniform or “weak” random key input. KDF is an important cryptographic primitive in practice and generate randomness for encryption, digital signatures, key exchange protocols, etc. HKDF [27] was introduced in 2010 by Krawczyk et al. as a

secure, generic cryptographic KDF and by now is the most deployed KDF. It follows the *extract-then-expand* paradigm [27]. For a given non-uniform or “weak” random key, HKDF extracts a relatively small but uniform and “highly” random string via an internal *randomness extractor*. Then, the extracted value is passed to a *randomness expander*: a pseudorandom function PRF with variable length output. The output is a uniform and “highly” random key-stream. HKDF is used in the Signal Protocol [34] for end-to-end encrypted messaging where it generates the message keys. There HKDF is used in conjunction with the triple Elliptic-curve Diffie-Hellman handshake [35] (X3DH) key agreement protocol. A number of applications of end-to-end encryption also make use of Signal and internally call the X3DH handshake and a secure (H)KDF to establish a key for later cryptographic use. Examples are the popular instant messaging apps WhatsApp [33], Facebook Messenger [2], Skype [30], Allo [32], Status [5], Secure Chat, Viber and Forsta. HKDF is a main component in the Noise Protocol Framework (NPF) [38], Message Layer Security (MLS) [3], and is used in TLS1.3.

HKDF is a key derivation function for processing general input-sources of some desired min-entropy. However, in many applications the sources are predefined, fixed and contain a nice algebraic structure, such as the triple Elliptic-curve Diffie-Hellman handshake [35] (X3DH) protocol in Signal. As a generic or non-source-specific KDF, HKDF meets basic security and performance goals, yet for concrete use cases, it might not be best optimized for: 1. source-specific (randomness) extraction; 2. performance on common platforms; 3. reliance on weaker assumptions than the random oracle for the SHA-256 hash function in HKDF. A dedicated KDF can efficiently leverage the structure of the input randomness source during the extraction phase to optimize the process. The possibility of constructing (input) source-specific extractor in a KDF under the *extract-then-expand* paradigm, was discussed in [27]. Concrete examples here are the works on deterministic randomness extraction from Diffie-Hellman schemes [14, 22].

To achieve secure and more efficient randomness expansion, a candidate building block is one that, unlike SHA-2, naturally expands its inputs. Cryptographically, such a function needs to come with some pseudorandom properties or behave as a fixed output length expanding pseudorandom function (PRF). Forkciphers [8] are such expanding functions, yet, they come with extra functionalities of inversion and reconstruction that are not necessary for randomness expansion and limit the PRF security to birthday-bound (in the

block size). The recently proposed PRF ButterKnife [7] is a fixed-length expanding PRF design and a natural KDF building block candidate. ButterKnife takes 256-bit inputs and produces 1024-bit outputs. It is based on AES and Deoxys-BC [25]. Its structure allows its implementation with AES native instructions (NI) on all AES-NI supporting processors. ButterKnife is proven secure both generically and analyzed cryptanalytically in [7]. ButterKnife’s security, in addition to its own security analysis, is further backed by the cryptanalysis results for AES-PRF [20, 37] and Deoxys-BC [16, 25, 29, 43].

*Expanding* primitives naturally fit better the concept of KDFs under the extract-then-expand paradigm, as opposed to the presently used compressing ones, e.g. SHA-2 in HKDF. The study of expanding PRFs for KDFs is cryptographically interesting in its own right and to the best of our knowledge, it is not done in the literature.

## 1.1 Our Results

In this work we propose Skye, a novel extract-then-expand KDF based on an expanding PRF. We provide a detailed security analysis of Skye and demonstrate empirically the efficiency advantages of Skye over HKDF, both directly and when used in Signal. We note that Skye applications are not limited to Signal and provide a discussion on other applications. Our contributions are:

**DETERMINISTIC EXTENSION FOR RANDOMNESS EXTRACTION.** We build a novel, generic and deterministic function  $\text{DExt}_f$  that aggregates the amount of extracted randomness for any randomness extractor  $f$  for multiple and independent input samples. We prove that the outputs of  $\text{DExt}_f$  are indistinguishable from random binary strings (see Sec. 5). Considering the randomness aggregation and extraction of the X3DH handshake of Signal input samples of multiple (three or four) DH shared secrets,  $\text{DExt}_f$  handles multiple independent samples and is hence well-suited for optimal randomness extraction. In this work,  $\text{DExt}_f$  is developed to be used in Skye for key derivation but we emphasize that this novel result is independent and quite interesting on its own.

**SECURE  $\text{DExt}_f$  INSTANTIATION FOR DH SAMPLES.** We provide a 128-bit secure simple and efficient *deterministic extractor* as an instantiation of  $\text{DExt}_f$ . This is achieved by choosing a DH source-specific extractor function to instantiate  $f$  in  $\text{DExt}_f$ . Applying the above-mentioned general result (in Sec. 5), we construct an extractor with a higher security margin from the msb/lbs (most/least significant bits based) extractor function for Diffie-Hellman (DH) schemes in [14]. We prove the security of our instantiation by combining the analysis of the generic  $\text{DExt}_f$  and the security of msb/lbs based extractor [14] (see Sec. 5.2).

**SECURE RANDOMNESS EXPANDER.** The outputs of the deterministic extractor are processed by a variable-output-length (VOL) PRF or randomness expander FExp. Two approaches to construct randomness expanders are presented in [27] – based on a counter or feedback encryption modes. HKDF uses as a randomness expander, a keyed feedback mode over the HMAC [28] pseudorandom function (PRF). To avoid the frequent rekeying of the feedback mode and the consequent speed reduction for the target 128-bit security level, we adopt a counter mode-like approach. FExp is motivated by the CTR\$ [39] encryption mode and it benefits from the expanding PRF function, as opposed to a block cipher. The FExp design achieves

both security and efficiency improvements over the CTR\$ [39] mode in the randomness expansion context. When compared to the CTR\$ mode, FExp accommodates larger and arbitrary (not necessarily random) inputs and provides full  $n = 128$ -bit security, as opposed to the birthday-bound  $n/2 = 64$ -bit, when instantiated with ButterKnife [7]. We give a security proof of FExp via the notion of indistinguishability from random binary strings, a notion equivalent to the IND\$ [39] notion (see Sec. 6).

**Skye: SECURE AND MORE EFFICIENT ALTERNATIVE TO HKDF IN SOURCE-SPECIFIC APPLICATIONS.** In Sec. 4, we propose the Skye scheme as a dedicated KDF following the extract-then-expand approach and using the two above-mentioned novel extractor and expander functions  $\text{DExt}_f$  and FExp, respectively.  $\text{Skye}[f, \text{PRF}_s]$  is based on two underlying primitives – a weak source-specific extractor  $f$  (processed under  $\text{DExt}_f$  extension) and an expanding PRF  $\text{PRF}_s$  (processed under FExp mode). We consider DH samples over Curve25519 as the input source to exemplify the performance and concrete security gain of Skye over HKDF and instantiate  $f$  with lsb. In Sec. 7, we show that when Skye is instantiated with ButterKnife, it achieves 128-bit CCS security [15, 27]. HKDF [27] is also shown secure under the CCS security notion and informally, CCS-security is defined as the indistinguishability of the KDF outputs from truly random strings under a chosen input attack. Here, we prove the security of Skye under a more practical (non random oracle) assumption than HKDF.

Our results of Theorem 5.1 and 7.1 show that similarly to HKDF, Skye can also be used with any type of input source with a well-defined weak extractor  $f$  (for that input source). More concretely, given an input source  $\Sigma$ , the KDF input set  $Z$  with elements sampled from  $\Sigma$ , an  $r$ -bit secure extractor  $f$  (w.r.t.  $\Sigma$ ) and an  $n$ -bit secure PRF  $\text{PRF}_s$ , the  $\text{Skye}[f, \text{PRF}_s]$  KDF provides  $\min\{n, c(r-1)\}$  bits of CCS security. Here  $c$  is a finite positive number that is directly proportional to the size of the input set  $Z$  and is inversely proportional to the desired extraction length (see Sec. 5.1 for details).

The security and performance of Skye for a given source depends on the underlying function  $f$  and hence it can be improved over HKDF for applications with known efficient extractors (as illustrated in this work with DH sources).

**SOFTWARE PERFORMANCE OF Skye.** We provide software performance comparison (in Sec. 9) of Skye with HKDF under their standard functionalities. Our results show that in isolation, Skye performs from 4x to 47x faster than HKDF in various settings depending on the availability of native instruction sets. We then consider Signal to further illustrate the performance benefits of Skye. We integrate Skye within the current implementation of Signal [1] and show that under various settings with and without available native instruction sets, Skye is able to provide 38-64% relative speedup in unidirectional messaging. In bidirectional messaging, where cost is dominated by the DH computation, the speedup depends on the number of messages sent at once. With a single message, Skye achieves 3%-11% relative speedup; after sending just 10 messages, the speedup rises to 12-36%, and with a further increase in the number of messages, it converges to the unidirectional speedup.

## 1.2 Related Works

**RANDOMNESS AMPLIFICATION.** The result of Maurer et al. [36, Corollary 2] can be applied to obtain similar security as  $\text{DExt}_f$  by iterating a set of independent input samples (with some weak extractor  $f$ ) and defining the  $\star$  operation is the XOR. Yet, our design approach is more general by allowing us to optimally extract (up to 100%) more randomness than that existing solution. For example, consider  $v$  input samples each providing  $w$  bits of extracted randomness and  $\lambda$  bits of security when passed through some weak extractor  $f$ . Now, if we need a higher security of  $c\lambda$  bits for some integer  $c$ , then the solution of [36] provides randomness only up to  $(v/c)w$  bits. Our solution  $\text{DExt}_f$  offers up to  $\lfloor (v-c)/\lceil c/2 \rceil \rfloor + 1$   $w$  bits of security which is optimal (see App. B) for various values of  $(v, c)$ . We generalize the result of Maurer et al. [36, Corollary 2] by offering a relatively simple proof that: 1. increases confidence in the existing analysis 2. generalizes and improves the bound (optimal) for reasonable/practical sample sizes.

**RELATED KDF NOTION.** In the formal analysis of the Double Ratchet protocol, Alwen et al. [4] formalized a KDF syntax called PRF-PRNG and its security notion called P. HKDF also satisfies this notion. For source-specific samples, our KDF Skye targets the stronger KDF notion of CCS-security, which implies the P-security when KDF is used as a PRF-PRNG (in the context of Signal). We refer the reader to App. D for more details on the PRF-PRNG syntax, how to convert a standard syntax KDF into a PRF-PRNG and a formal claim with full proof of CCS to P security reduction.

Alwen et al. [4] provide a standard model proof for the Double Ratchet and Signal Protocol, however, their analysis is not RO-free when the KDF is instantiated with HKDF or alternatives with idealized assumptions (see App. D). This leads us to naturally consider the possibility of designing a KDF for applications like (but not limited to) the Signal Protocol that is RO-free and is equally or more efficient than the current HKDF solution. The Skye KDF follows the standard KDF syntax [27] and is proven to be CCS [15, 27]-secure in the standard model.

## 2 NOTATIONS

**STRINGS.** All strings are binary strings. The set of all strings of length  $n$  (for a positive integer  $n$ ) is denoted by  $\{0, 1\}^n$  and the set of all strings of all possible lengths is denoted by  $\{0, 1\}^*$ . We denote by  $\text{Func}(m, n)$  the set of all functions with domain  $\{0, 1\}^m$  and range  $\{0, 1\}^n$ . For a string  $X$  of  $\ell$  bits, we let  $X[i]$  denote the  $i^{\text{th}}$  bit of  $X$  for  $i = 0, \dots, \ell - 1$  (starting from the left) and  $X[i \dots j] = X[i] \parallel X[i+1] \parallel \dots \parallel X[j]$  for  $0 \leq i < j < \ell$ . We let  $\text{msb}_\ell(X) = X[0 \dots (\ell - 1)]$  denote the  $\ell$  leftmost (most significant) bits of  $X$  and  $\text{lsb}_r(X) = X[(|X| - r) \dots (|X| - 1)]$  the  $r$  rightmost (least significant) bits of  $X$ , such that  $X = \text{msb}_\chi(X) \parallel \text{lsb}_{|X| - \chi}(X)$  for any  $0 \leq \chi \leq |X|$ . Given an  $xn$ -bit string  $X$ , we let  $X_1, \dots, X_x \stackrel{n}{\leftarrow} X$  denote a partitioning of  $X$  into  $n$ -bit blocks, such that  $|X_i| = n$  for  $i = 1, \dots, x$ .

**MISCELLANEOUS.** The symbol  $\perp$  denotes an error signal, or an undefined value. We denote by  $X \stackrel{\$}{\leftarrow} \mathcal{X}$ , a sampling of an element  $X$  from a finite set  $\mathcal{X}$  following the uniform distribution. We use

lexicographic comparison of tuples of integers; i.e.  $(i', j') < (i, j)$  iff  $i' < i$  or  $i' = i$  and  $j' < j$ . We denote the set of natural numbers by  $\mathbb{N}$ . For a matrix  $M$ , we use  $M[ij]$  to denote the  $ij^{\text{th}}$  entry of  $M$ . For equations  $E_i$ s with  $1 \leq i \leq a$  (for some  $a \in \mathbb{N}$ ), we use the indexed set  $S = \{E_i | 1 \leq i \leq a\}$  to denote the system of equations  $E_i$ s.

## 3 KEY DERIVATION FUNCTION

KDFs derive one or more *cryptographically secure* secret keys (of any fixed length) from a source of initial key material (IKM) that can contain some good amount of randomness but is not distributed uniformly. The notion of cryptographically secure keys is usually associated with pseudorandom keys, i.e. keys that are computationally indistinguishable from a uniformly random string of the same length. Formally, a source of IKM can be defined as follows.

**DEFINITION 3.1 (SOURCE OF IKM [27]).** A source of initial keying material (or simply source)  $\Sigma$  is a two-valued probability distribution  $(Z, C_Z)$  generated by an efficient probabilistic algorithm (we will refer to both the probability distribution as well as the generating algorithm by  $\Sigma$ ).

Here,  $Z$  values are used to denote the random secret IKM samples with required entropy (e.g., DH shared secrets) whereas  $C_Z$  provides the public info associated to  $Z$  (e.g., DH public keys and group parameters). See Sec. 7 for more details.

A KDF function tackles the case when the initial key material is not pseudorandom or uniformly random, e.g. the initial key material is obtained from a weak process that uses renewable sources of randomness, a weak random number generator, random sampling over a group or DH values computed in a (key-exchange) protocol. In these settings a KDF function is constructed using the so-called *extract-then-expand* paradigm.

### 3.1 Extract-then-Expand Paradigm

A KDF following the extract-then-expand paradigm has two components: 1. a randomness extractor  $\text{Ext}$  that extracts a fixed-length pseudorandom key  $K$  from an “imperfect” source of initial key material; 2. randomness expander  $\text{Exp}$  that expands  $K$  to a variable-length output. The latter is usually built using a regular pseudorandom function (PRF) with output extension via counter or feedback encryption modes [27, 39].  $\text{Ext}$  produces “close-to-random” outputs (in a computational or statistical sense) from an input that is sampled from the corresponding source key distribution. The extraction may have an additional non-secret input or salt value, that is either randomized or constant. When it is constant, the extractor and the corresponding KDF are called *deterministic extractor* and *deterministic KDF*, respectively. The expander  $\text{Exp}$  uses  $K$  to produce cryptographic keys  $K_{kdf}$  of a specified length.  $\text{Exp}$  takes the output length parameter as one of the inputs.

A KDF scheme is defined over the inputs: a source key  $Z$ , the extractor salt  $\text{salt}$  (null or constant), the length  $\ell$  of key bits to be output by KDF, and a context variable or auxiliary info string  $\gamma$  (may be null). The latter string includes key-related information that is uniquely (and cryptographically) bound to the produced output. For example, includes information about the application or protocol calling the KDF, session-specific information like nonces,

time, session identifiers, etc. The KDF evaluation is defined as:  $K = \text{Ext}(\text{salt}; Z)$ ,  $K_{kdf} = \text{Exp}(K; \gamma; \ell)$ .

**HKDF KEY DERIVATION FUNCTION.** In the expansion phase the output is  $K_{kdf} = K_1 || K_2 || \dots || K_t[1 \dots \ell]$  where  $t = \lceil \ell / n_{\text{HMAC}} \rceil \in \mathbb{N}$ ,  $\ell$  is the desired output length and  $n_{\text{HMAC}}$  is the HMAC output size.  $K_i$  are obtained as:

$$K_1 = \text{HMAC}(K, \gamma || 0), K_{i+1} = \text{HMAC}(K, K_i || \gamma || i), \quad 1 \leq i < t,$$

where  $\gamma$  denotes the context variable.

### 3.2 Deterministic Multi-sample KDF

A KDF is *multi-sample* if it takes multiple independent samples from its input source and combines those for relatively large and highly random output with the same or increased security than its single-sample variant. Formally, a deterministic multi-sample KDF  $\Pi : (\{0, 1\}^*)^v \times \Gamma \times \mathbb{N} \rightarrow \{0, 1\}^*$  takes three arguments; a set  $Z$  of  $v$  values  $Z_i$ s (binary strings) sampled from a source of keying material (defined in Def. 3.1), a context variable or auxiliary info  $\gamma$  (optional, i.e., null string or a constant) and the desired output length  $\ell$ . The function returns an  $\ell$ -bit binary string  $K_{kdf}$ . HKDF [27] is a deterministic multi-sample KDF scheme when its salt value is set to a constant.

## 4 Skye: AN EXPANDING PRF BASED KDF

In this section, we define a new KDF called Skye. Skye is based on an expanding pseudorandom function  $\text{PRF}_s$  where  $\text{PRF}_s : \{0, 1\}^k \times \{0, 1\}^{2n} \rightarrow \{0, 1\}^{sn}$  is a symmetric primitive that transforms a fixed length (2n-bit) input  $X$  into a larger fixed-length ( $sn$ -bit with  $s \geq 2$ ) output  $Y$  via a secret key  $K \in \mathcal{K}$  of  $k$  bits.

**DEFINITION 4.1 (PRF SECURITY).** For  $\text{PRF}_s : \{0, 1\}^k \times \{0, 1\}^{2n} \rightarrow \{0, 1\}^{sn}$ , let  $\mathcal{A}$  be an adversary whose goal is to distinguish  $\text{PRF}_s(K, \cdot)$  and a uniform random function  $R(\cdot) : \{0, 1\}^{2n} \rightarrow \{0, 1\}^{sn}$  by their oracle access. The prf-security of  $\text{PRF}_s$  is defined as:

$$\text{Adv}_{\text{PRF}_s}^{\text{prf}}(\mathcal{A}) = |\Pr[K \xleftarrow{\$} \mathcal{K} : \mathcal{A}^{\text{PRF}_s(K, \cdot)} \Rightarrow 1] - \Pr[R \xleftarrow{\$} \text{Func}(2n, sn) : \mathcal{A}^{R(\cdot)} \Rightarrow 1]|.$$

**DESCRIPTION OF Skye.** Skye is a deterministic multi-sample KDF that follows the extract-then-expand approach by using a deterministic source-specific extractor  $\text{DExt}_f$  and an expanding prf  $\text{PRF}_s$ -based randomness expander  $\text{FExp}$ , respectively.

More specifically, for inputs; a set of  $v$  binary strings  $Z$ , a context variable  $\gamma$  (from some finite set  $\Gamma$ ) and an output length  $\ell \in \mathbb{N}$ ,  $\text{Skye} : (\{0, 1\}^*)^v \times \Gamma \times \mathbb{N} \rightarrow \{0, 1\}^*$  is defined as

$$\text{Skye}(Z, \gamma, \ell) = \text{FExp}(K_{\text{ext}}, \gamma, \ell) = K_{\text{exp}}$$

where  $K_{\text{ext}} = \text{DExt}_f(Z)$  and  $|K_{\text{exp}}| = \ell$ .

In the following sections, we define the two components  $\text{DExt}_f$  and  $\text{FExp}$  in full detail using their corresponding underlying functions; the deterministic extractor  $f$  and the expanding PRF  $\text{PRF}_s$ , respectively. With these components defined, any instantiation of  $\text{Skye}[f, \text{PRF}_s]$  can be solely described by the two primitive functions  $f$  and  $\text{PRF}_s$ .

For real instantiations (see Sec. 8), we use ButterKnife [7] as the  $\text{PRF}_s$  in Skye with  $n = 128$  bits and  $s = 8$ . ButterKnife is well-suited

due to its AES internal structure and as a result benefiting from the AES native instructions (NI) on all supporting processors. Additionally, it has a very large expansion of  $s = 8$  for performance gains and efficient key-scheduling [7] that saves the cost of frequent key scheduling. ButterKnife comes with dedicated cryptanalysis and a generic proof of security [7] and is backed by further cryptanalytic results for AES-PRF [20, 37] and Deoxys-BC [16, 25, 29, 43].

In contrast to  $\text{PRF}_s$ , instantiating  $f$  is not that simple as it depends on the KDF application and the input source of extraction. In other words, a secure extractor for some input source may not work or be secure for others. A large set of KDF applications as mentioned in Sec. 1 deal with multi-DH sources such as X3DH based applications or Signal-like protocols. Therefore, to illustrate the performance and security advantages of Skye in these applications, we consider multi-DH as the exemplary input source and instantiate  $f$  with the least significant bits function  $\text{lsb}$  (analyzed for multi-DH sources in Sec. 5.2).

For brevity, in the remainder of the paper, we drop  $f$  from the global inputs whenever it is fixed to  $\text{lsb}$ , i.e. when we are working with multi-DH sources. We note that for any other source (i.e. applications beyond multi-DH sources) with a suitable weak extractor  $f_0$ , Skye can be simply defined by reinstantiating  $f$  with  $f_0$  instead of  $\text{lsb}$ .

## 5 RANDOMNESS EXTRACTION PHASE

We first recall some preliminary definitions that are needed for the upcoming results.

**DEFINITION 5.1 (STATISTICAL DISTANCE).** Let  $X$  and  $Y$  be two random variables taking values from a finite set  $\mathcal{X}$ . The statistical distance between  $X$  and  $Y$  is the value of the following expression:

$$\text{SD}(X, Y) = \frac{1}{2} \sum_{x \in \mathcal{X}} |\Pr[X = x] - \Pr[Y = x]|.$$

### ELLIPTIC CURVE DECISIONAL DIFFIE-HELLMAN (ECDDH) PROBLEM [11].

Let  $G$  be a public cyclic subgroup (of size a prime  $q$ ) of an elliptic curve group  $\mathcal{E}(\mathbb{F}_p)$  over the finite field  $\mathbb{F}_p$  of size prime  $p$ . Let  $P$  be a randomly chosen generator of  $G$ , then for some secret and randomly chosen  $a, c \in \mathbb{Z}_q^*$  and a randomly chosen  $Q \in G$ , it is hard to determine whether  $Q = acP$  or  $Q \neq acP$  when provided with  $P, aP, cP$  and  $Q$ . Formally, we define the advantage of an adversary  $\mathcal{A}$  against solving ECDDH in  $G$  as

$$\text{Adv}_G^{\text{ecddh}}(\mathcal{A}) = |\Pr[\mathcal{A}(P, aP, cP, acP) = 1] - \Pr[\mathcal{A}(P, aP, cP, Q) = 1]|$$

where the probability is over the random choices of  $a, c \in \mathbb{Z}_q^*$ ,  $P, Q \in G$  and the internal randomness of  $\mathcal{A}$ .

The ECDDH assumption (with respect to  $G$ ) states that for all efficient adversaries, the value of advantage  $\text{Adv}_G^{\text{ecddh}}(\mathcal{A})$  is reasonably small.

In the remaining section, we now show how to build a strong extractor from a weak extractor and provide its security analysis. Formally, an extractor is defined as follows:

**DEFINITION 5.2 ((X, ε)-EXTRACTOR [27]).** Let  $X$  be a random variable that takes values in  $\{0, 1\}^x$  and  $U_w$  denote a random variable uniformly distributed in  $\{0, 1\}^w$  for some positive integers  $x$  and  $w$ ,

where  $U_w$  and  $X$  are independent. We say that for some non-negative integer  $d$ , the function  $f : \{0, 1\}^x \times \{0, 1\}^d \rightarrow \{0, 1\}^w$  is an  $(X, \epsilon)$ -extractor if

$$\text{SD}((f(X, \text{salt}), \text{salt}), (U_w, \text{salt})) \leq \epsilon$$

where the salt is a  $d$ -bit value which is made public upon sampling.

We note that the salt value can be a null string or constant (in case of deterministic extractors). Further, since the salt sampling is independent to  $f, X$  and  $U_w$ , for simplicity, we slightly abuse the notation and omit both salt terms from the second arguments in the expression of Def. 5.2. In other words, the expression becomes

$$\text{SD}(f(X, \text{salt}), U_w) \leq \epsilon.$$

### 5.1 DExt: A Generic and Deterministic Extension Towards Improved Security of any Extractor

In this section, we describe how to build a stronger extractor  $\text{DExt}_f$  which can more securely extract the same or larger amount of randomness from a relatively weak extractor  $f$  given multiple independent samples. Our construction uses only simple and cheap operators like concatenation and XOR to avoid additional computational cost. Unlike the prior works based on extraction from multiple sources [17, 18, 21, 26, 41], our construction requires independence only within the samples where the source of these samples can still remain the same. An example here is multiple independent DH handshakes defined over the same source group.

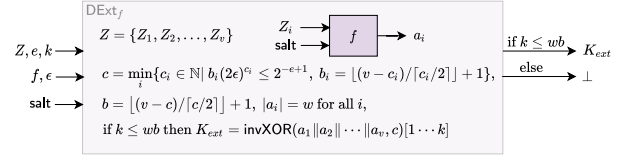
**DEFINITION OF  $\text{DExt}_f$ .** We provide a definition of the extended extractor  $\text{DExt}_f$  in Fig. 1 as a function that takes as inputs:

- a set  $Z = \{Z_1, Z_2, \dots, Z_v\}$  of  $v$  many  $z$ -bit (when represented in binary) independent samples chosen from some public finite sets  $S_1, S_2, \dots, S_v$ , respectively (all of these sets could be same) for some positive integer  $z$ ;
- the desired output length  $k$ ;
- a positive integer  $e$  (security parameter);
- a pair  $(f, \epsilon)$  with  $f : \{0, 1\}^z \times \{0, 1\}^d \rightarrow \{0, 1\}^w$  being a  $(U_{Z_i}, \epsilon)$ -extractor for some  $\epsilon > 0$ , positive integers  $d$  and  $w$  and all  $1 \leq i \leq v$ ;
- a  $d$ -bit salt value salt (if any).

$\text{DExt}_f$  then uses the  $(U_{Z_i}, \epsilon)$ -extractor (for all  $1 \leq i \leq v$ )  $f$  and the function  $\text{invXOR}$  (defined in Def. 5.3) and outputs a  $k$ -bit string  $K_{\text{ext}}$  if  $k \leq wb$  or  $\perp$  otherwise. Here,  $U_{Z_i}$  denotes a random variable distributed according to the sampling of  $Z_i$ s from the set  $S_i$  for all  $i$ s, respectively,  $b = \lfloor \frac{v-c}{\lceil c/2 \rceil} \rfloor + 1$  and  $c = \min_i \{c_i \in \mathbb{N} \mid b_i (2\epsilon)^{c_i} \leq 2^{-e+1}, b_i = \lfloor \frac{v-c_i}{\lceil c_i/2 \rceil} \rfloor + 1\}$ . We note that the parameters  $c, e, b$  and the function  $f$  are important to define  $\text{DExt}_f$ , its design optimality (Theorem B.2) and security (Theorem 5.1 and its following corollary).

**DEFINITION 5.3 (XOR-BASED INVOLVEMENT).** We define  $\text{invXOR}$  as a function that maps a  $wv$ -bit binary string  $a_1 \| a_2 \| \dots \| a_v$  where all  $a_i$ s are  $w$ -bit binary strings (equivalently elements of  $\mathbb{F}_{2^w}$  for some integer  $w \geq 0$ ) and an integer  $c < v$  to an  $wb$ -bit binary string as follows:

$$\text{invXOR}(a_1 \| a_2 \| \dots \| a_v, c) = \oplus_{i=1}^c a_i \| \oplus_{i=1+\lceil c/2 \rceil}^{i+c+\lceil c/2 \rceil} a_i \| \dots \| \oplus_{i=1+(b-1)\lceil c/2 \rceil}^{i+c+(b-1)\lceil c/2 \rceil} a_i$$



**Figure 1:  $\text{DExt}_f$ : A  $(U_Z, b(2\epsilon)^c/2)$ -extractor (with upto  $wb$  bits of extraction) parameterized by a relatively small and weak  $(U_{Z_i}, \epsilon)$ -extractor  $f$  (with upto  $w$  bits of extraction). Here  $U_{Z_i}$  and  $U_Z$  denote random variables distributed according to the sampling of  $Z_i$ s from  $S_i$ s for all  $i$ s and of corresponding  $Z$  from the Cartesian product of all  $S_i$ s, respectively.**

where  $b = \lfloor (v - c) / \lceil c/2 \rceil + 1 \rfloor$ .

We refer the reader to App. B for a detailed analysis and discussion on the (sub-)optimality of the  $\text{DExt}_f$  construction. This analysis shows that although finding an optimal extension for randomness extraction over arbitrary values of  $v, c$  is hard, it is possible to construct one for  $v < 6$  (the relevant case in practice) and  $\text{DExt}_f$  is one example of such extensions.

With all necessary definitions in place, we are now ready to give the formal security statement of  $\text{DExt}_f$  in Theorem 5.1.

**THEOREM 5.1.** Let  $Z_1, Z_2, \dots, Z_v$  be  $v$  independent  $z$ -bit binary strings that are chosen from the public finite sets  $S_1, S_2, \dots, S_v$ , respectively for some positive integer  $z$ . Let  $f$  be a public function with a range in  $\{0, 1\}^w$  for some positive integer  $w$ . Let  $k$  and  $c$  be two positive integers, such that  $c \leq v$  and  $k \leq wb$  where  $b = \lfloor (v - c) / \lceil c/2 \rceil + 1 \rfloor$ . If there exists an  $\epsilon > 0$  such that  $f$  is a  $(U_{Z_i}, \epsilon)$ -extractor for all  $i$  then, we have

$$\text{SD}(U_D, U_k) \leq \frac{b}{2} (2\epsilon)^c.$$

Here  $U_{Z_i}$  is a random variable distributed according to the sampling of  $Z_i$ s in the set  $S_i$ .  $U_k$  is a random variable uniformly distributed over  $\{0, 1\}^k$ .  $U_D$  is a random variable distributed according to the outputs of  $\text{DExt}_f$  (of size  $k$  bits) when provided with  $Z_1, Z_2, \dots, Z_v, k, c$  and salt (if any) as inputs.

Increasing the value of  $c$  for a fixed distribution of  $Z_i$ s will increase the security margin of the final outputs linearly, but will also hyperbolically decrease the maximum possible length of the final outputs. This trade-off is important and explains why it is good to leave  $c$  a free variable in the theorem (this can later be defined according to the requirements of an application).

**PROOF OF THEOREM 5.1.** We are given  $v$  independent  $z$ -bit elements  $Z_1, \dots, Z_v$  (when represented as binary strings) that are chosen from some public finite sets  $S_1, S_2, \dots, S_v$ , respectively for some positive integer  $z$  and a public function  $f$  (that may or may not require a random salt value for its evaluation) with a range in  $\{0, 1\}^w$  for some positive integer  $w$ . Let us now consider  $U_w$  as a random variable uniformly distributed over  $\{0, 1\}^w$ . For some  $\epsilon > 0$ , we have  $\text{SD}(f(U_{Z_i}, \text{salt}), U_w) \leq \epsilon$  for all  $i$ . Since the value of salt is sampled once and used for all values of  $Z_i$ , for simplicity, we denote  $f(U_{Z_i}, \text{salt})$  by  $f(U_{Z_i})$  in the rest of the proof. Let  $c, k$  be two positive integers and  $b = \lfloor (v - c) / \lceil c/2 \rceil + 1 \rfloor$  such that  $w = \lceil kb \rceil$ .

Now, from the definition of SD and  $\epsilon$  we have

$$\text{SD}(f(U_{Z_i}), U_w) = \frac{1}{2} \sum_{x \in U_w} \left| \Pr[f(U_{Z_i}) = x] - \frac{1}{2^w} \right| \leq \epsilon. \quad (1)$$

Also, since an element of  $U_D$  of size  $k$  that corresponds to  $Z_1, \dots, Z_v$  can be equivalently defined as  $\text{invXOR}(f(Z_1) \| f(Z_2) \| \dots \| f(Z_v), c) [1 \dots k]$ , we have for  $k' = b \lceil k/b \rceil$ ,

$$\begin{aligned} \text{SD}(U_D, U_k) &\leq \text{SD}(\text{invXOR}(f(U_{Z_1}) \| f(U_{Z_2}) \| \dots \\ &\quad \dots \| f(U_{Z_v}), c), U_{k'}) \\ &= \text{SD}(\oplus_{i=1}^c f(U_{Z_i}) \| \oplus_{i=1+\lceil c/2 \rceil}^{i=c+\lceil c/2 \rceil} f(U_{Z_i}) \| \dots \\ &\quad \dots \| \oplus_{i=1+(b-1)\lceil c/2 \rceil}^{i=c+(b-1)\lceil c/2 \rceil} f(U_{Z_i}), U_{k'}) \end{aligned}$$

We denote  $\oplus_{i=1+(j-1)\lceil c/2 \rceil}^{i=c+(j-1)\lceil c/2 \rceil} f(U_{Z_i})$  by  $U_{f_j}$ . Hence, we have

$$\begin{aligned} \text{SD}(U_D, U_k) &\leq \text{SD}(U_{f_1} \| U_{f_2} \| \dots \| U_{f_b}, U_{k'}) \\ &= \frac{1}{2} \sum_{x \in U_{k'}} |\Pr[U_{f_1} \| U_{f_2} \| \dots \| U_{f_b} = x] - \Pr[U_{k'} = x]| \\ &= \frac{1}{2} \sum_{\substack{x \in U_{k'} \\ x_1, \dots, x_b \xleftarrow{w} x}} \left| \prod_{j=1}^b \Pr[U_{f_j} = x_j] - \frac{1}{2^{k'}} \right|. \end{aligned}$$

The last equality holds due to the fact that for each  $U_{f_j}$  the corresponding subset of  $Z_i$ s contains at least one new/fresh independent element from the corresponding main set  $Z$  (note that this is true for all positive values of  $c$ ). Hence all  $U_{f_j}$ s can be considered independent from each other. Now, with some basic algebra we can show that for all  $y_j$ s

$$\prod_{j=1}^b y_j - \frac{1}{2^{k'}} = \sum_{j=1}^b \left( \left( \frac{1}{2^w} \right)^{j-1} \left( y_j - \frac{1}{2^w} \right) \prod_{j'=j+1}^b y_{j'} \right).$$

Hence, we have  $\text{SD}(U_D, U_k)$

$$\begin{aligned} &\leq \frac{1}{2} \sum_{\substack{x \in U_{k'} \\ x_1, \dots, x_b \xleftarrow{w} x}} \sum_{j=1}^b \left| \left( \frac{1}{2^w} \right)^{j-1} \left( \Pr[U_{f_j} = x_j] - \frac{1}{2^w} \right) \prod_{j'=j+1}^b \Pr[U_{f_{j'}} = x_{j'}] \right| \\ &= \sum_{j=1}^b \frac{1}{2^{wj-w+1}} \sum_{\substack{x \in U_{k'} \\ x_1, \dots, x_b \xleftarrow{w} x}} \left| \left( \Pr[U_{f_j} = x_j] - \frac{1}{2^w} \right) \prod_{j'=j+1}^b \Pr[U_{f_{j'}} = x_{j'}] \right| \\ &= \sum_{j=1}^b \frac{1}{2^{wj-w+1}} \left( \sum_{x_j \in U_w} \left| \Pr[U_{f_j} = x_j] - \frac{1}{2^w} \right| \right) \\ &\quad \left( \prod_{j'=j+1}^b \sum_{x_{j'} \in U_w} \Pr[U_{f_{j'}} = x_{j'}] \right) \left( \prod_{j'=1}^{j-1} \sum_{x_{j'} \in U_w} 1 \right) \\ &= \sum_{j=1}^b \text{SD}(U_{f_j}, U_w). \end{aligned} \quad (2)$$

**CLAIM 1 (1).** For  $\epsilon > 0$  defined as above we have  $\text{SD}(U_{f_j}, U_w) \leq (2\epsilon)^c/2$  for all  $1 \leq j \leq b$ .

Let us assume for a moment that Claim (1) holds then combining this result with Eqn. 2 gives us the result of Theorem 5.1 and completes its proof. We now prove the result of Claim (1). Let us

recall that  $U_{f_j} = \oplus_{i=1+(j-1)\lceil c/2 \rceil}^{i=c+(j-1)\lceil c/2 \rceil} f(U_{Z_i})$  for all  $1 \leq j \leq b$ . Again, for simplicity, we denote  $1 + (j-1)\lceil c/2 \rceil$  and  $c + (j-1)\lceil c/2 \rceil$  by  $\alpha$  and  $\beta$ , respectively. This implies

$$\begin{aligned} \text{SD}(U_{f_j}, U_w) &= \text{SD}(\oplus_{i=\alpha}^{\beta} f(U_{Z_i}), U_w) \\ &= \frac{1}{2} \sum_{x_j \in U_w} \left| \Pr[\oplus_{i=\alpha}^{\beta} f(U_{Z_i}) = x_j] - \frac{1}{2^w} \right| \\ &= \frac{1}{2} \sum_{x_j \in U_w} \left| \sum_{\substack{\oplus_{i=\alpha}^{\beta} x_{ij} = x_j \\ x_{ij} \in f(U_{Z_i})}} \Pr[\wedge_{i=\alpha}^{\beta} (f(U_{Z_i}) = x_{ij})] - \frac{1}{2^w} \right|. \end{aligned}$$

As we know, the system of linear equations that corresponds to the outputs of  $\text{invXOR}$  (i.e. a system with equations defined by the concatenated blocks of an  $\text{invXOR}$  output) has  $\text{mdi} = c - 1$ , which means every variable  $f(U_{Z_i})$  in the equation  $E_j := \oplus_{i=\alpha}^{\beta} f(U_{Z_i}) = x_j$  has at least  $c - 1$  degree of involvement. Now, since for all  $j$ s,  $E_j$  contains exactly  $\beta - \alpha = c - 1$  variables, we have for all  $E_j$ s and  $x_{ij}$  as defined above

$$\Pr[\wedge_{i=\alpha}^{\beta} (f(U_{Z_i}) = x_{ij})] = \prod_{i=\alpha}^{\beta} \Pr[f(U_{Z_i}) = x_{ij}].$$

Hence, we get  $\text{SD}(U_{f_j}, U_w)$

$$\begin{aligned} &= \frac{1}{2} \sum_{x_j \in U_w} \left| \sum_{\substack{\oplus_{i=\alpha}^{\beta} x_{ij} = x_j \\ x_{ij} \in f(U_{Z_i})}} \prod_{i=\alpha}^{\beta} \Pr[f(U_{Z_i}) = x_{ij}] - \frac{1}{2^w} \right| \\ &= \frac{1}{2} \sum_{x_j \in U_w} \left| \sum_{\substack{\oplus_{i=\alpha}^{\beta} x_{ij} = x_j \\ x_{ij} \in f(U_{Z_i})}} \prod_{i=\alpha}^{\beta} \Pr[f(U_{Z_i}) = x_{ij}] - \sum_{\substack{\oplus_{i=\alpha}^{\beta} x_{ij} = x_j \\ x_{ij} \in U_w}} \prod_{i=\alpha}^{\beta} \frac{1}{2^w} \right| \\ &\leq \frac{1}{2} \sum_{x_j \in U_w} \left| \sum_{\substack{\oplus_{i=\alpha}^{\beta} x_{ij} = x_j \\ x_{ij} \in U_w}} \prod_{i=\alpha}^{\beta} \Pr[f(U_{Z_i}) = x_{ij}] - \sum_{\substack{\oplus_{i=\alpha}^{\beta} x_{ij} = x_j \\ x_{ij} \in U_w}} \prod_{i=\alpha}^{\beta} \frac{1}{2^w} \right| \\ &\leq \frac{1}{2} \sum_{x_j \in U_w} \sum_{\substack{\oplus_{i=\alpha}^{\beta} x_{ij} = x_j \\ x_{ij} \in U_w}} \prod_{i=\alpha}^{\beta} \left| \Pr[f(U_{Z_i}) = x_{ij}] - \frac{1}{2^w} \right| \\ &= \frac{1}{2} \sum_{x_{ij} \in U_w} \prod_{i=\alpha}^{\beta} \left| \Pr[f(U_{Z_i}) = x_{ij}] - \frac{1}{2^w} \right| \\ &= \frac{1}{2} \prod_{i=\alpha}^{\beta} 2 \cdot \text{SD}(f(U_{Z_i}), U_w) \leq \frac{1}{2} (2\epsilon)^c. \end{aligned}$$

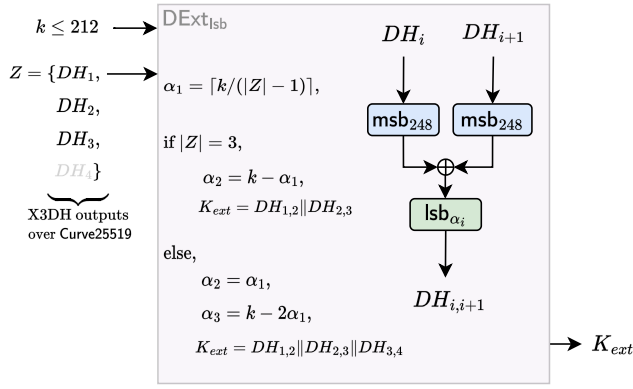
Here the first inequality holds because  $|f(U_{Z_i})| \leq |U_w| = 2^w$  for all  $1 \leq i \leq v$  and the last inequality follows from Eqn. 1.  $\square$

## 5.2 DExt<sub>f</sub> Instantiation based on msb/lsh Function

In this section, we show that for applications with multi-DH input sources such as the Signal Protocol where the KDF source of



inputs is a subgroup of an elliptic curve group and the inputs are DH elements over it i.e. the input distribution is computationally indistinguishable from the uniform distribution over the source, one can instantiate the underlying function  $f$  of  $\text{DExt}_f$  by the functions  $\text{msb}$  or  $\text{lsb}$ . More concretely, we recall one of the two main theorems from the work of Chevalier et al. [14] on the security of  $\text{lsb}_k$  function as a deterministic extractor (this function does not require any additional salt during its evaluation) and combine it with Theorem 5.1 to amplify the extracted output's size and security.



**Figure 2:  $\text{DExt}_{\text{lsb}}$ : A deterministic extractor that can extract up to 212 bits of randomness with 128 bits of security from 3 or 4 random, independent and fresh DH shared secrets that are X3DH outputs over Curve25519. Here  $\text{msb}$  and  $\text{lsb}$  are the most and least significant bits functions, respectively, as defined in Sec. 2.**

**THEOREM 5.2 (msb/lsb EXTRACTION, THEOREM 14, [14]).** *Let  $p$  be an  $\ell_p$ -bit prime,  $G$  a subgroup of  $\mathcal{E}(\mathbb{F}_p)$  of cardinality  $q$  generated by  $P_0$ ,  $q$  being an  $\ell_q$ -bit prime,  $U_G$  and  $U_k$  be two random variables uniformly distributed in  $G$  and  $\{0, 1\}^k$ , respectively for some positive integer  $k$ . Then we have*

$$\text{SD}(\text{lsb}_k(U_G), U_k) \leq 2^{(k+\ell_p+\log_2 \ell_p)/2+3-\ell_q}.$$

By combining the results from Theorem 5.1 and 5.2 we obtain the following corollary:

**COROLLARY 5.2.1.** *Let  $p$  be an  $\ell_p$ -bit prime,  $G$  a subgroup of  $\mathcal{E}(\mathbb{F}_p)$  of cardinality  $q$  generated by  $P_0$ ,  $q$  being an  $\ell_q$ -bit prime,  $U_D$  a random variable distributed according to the outputs of  $\text{DExt}_{\text{lsb}}$  for  $v$  many chosen uniform and independent samples from  $G$  with two positive integers  $k$  and  $c(< v)$  and  $b = \lfloor (v-c)/\lceil c/2 \rceil \rfloor + 1$ . Let  $U_k$  be a random variable uniformly distributed in  $\{0, 1\}^k$ . We have  $\text{SD}(U_D, U_k) \leq 2^{-e}$  for some positive integer  $e$  if*

$$2\ell_q - \ell_p - \log_2 \ell_p - 6 \geq \frac{2(e + \log_2 b - 1)}{c} + \left\lceil \frac{k}{b} \right\rceil.$$

In the concrete settings of the Signal Protocol, we have the source group  $G$  defined as the cyclic subgroup of Curve25519 [10] using the base point  $x = 9$  (one of the NIST standards for ECC [13] targeting 128-bit security) with  $\ell_p = 256$  bits and  $253 \geq \ell_q \geq 252$ . Further, under the ECDDH [11] assumption (see Sec. 5), the X3DH [35]

handshake over  $G$  provides at least 3 uniform and independent group elements as DH shared secrets. Hence, from Corollary 5.2.1, we have that  $\text{DExt}$  can output upto  $k = 212$  bits of randomness with security of  $e = 128$  bits when provided with fresh, independent and random X3DH outputs with the least significant byte (a.k.a. the clamped [31] byte that contains three fixed bits as 0s) dropped.

Note that with the results from [14] (Theorem 5.2) one can achieve the same amount of randomness, i.e.  $k \approx 212$  bits by concatenating the outputs of  $f = \text{lsb}$ , but with at most security of  $e = 82$  bits when provided with at least  $v = 3$  uniform and independent samples from the group  $G$ .

Let us consider the following example to clarify a different implication of this result. For a fixed sample size  $v = 10$ , source settings  $p = 256$  and  $q \geq 252$ , and a security parameter  $e = 80$ , the input size that can be extracted in the final concatenated output (computed as concatenation of the outputs of  $f = \text{lsb}$ ) is only  $\approx 26\%$  (according to the results of [14]). On the other hand, the same is improved to 52% (according to our results) when evaluated with  $\text{DExt}_{\text{lsb}}$  under the same settings.

To summarize, for a given X3DH output with 3 (resp. 4) random, independent and fresh DH samples (defined after dropping the fixed/clamped byte) as  $\text{DH1}||\text{DH2}||\text{DH3}$  (resp.  $\text{DH1}||\text{DH2}||\text{DH3}||\text{DH4}$ ) over Curve25519, we have shown that under the ECDDH assumption the string  $\text{lsb}_{\lceil k/2 \rceil}(\text{DH1} \oplus \text{DH2})||\text{lsb}_{k-\lceil k/2 \rceil}(\text{DH2} \oplus \text{DH3})$  (resp.  $\text{lsb}_{\lceil k/3 \rceil}(\text{DH1} \oplus \text{DH2})||\text{lsb}_{\lceil k/3 \rceil}(\text{DH2} \oplus \text{DH3})||\text{lsb}_{k-2\lceil k/3 \rceil}(\text{DH3} \oplus \text{DH4}))$  is indistinguishable from a  $k$ -bit uniform random string with a security of at least 128 bits when  $k \leq 212$ . A simplified diagram of  $\text{DExt}_f$  with  $f$  defined as the  $\text{lsb}$  function for Signal-like applications is provided in Fig. 2.

## 6 RANDOMNESS EXPANSION PHASE

Randomness expansion or key-stream generation is a process to generate a large amount of random bits from relatively smaller random and secret keys. Formally, a randomness expander scheme  $\Pi : \mathcal{K} \times \Gamma \times \mathbb{N} \rightarrow \{0, 1\}^*$  takes a  $k$ -bit key  $K \in \mathcal{K}$ , the desired output length  $\ell \in \mathbb{N}$  and an additional arbitrary but fixed length binary string  $\gamma \in \Gamma$  as inputs and returns an  $\ell$ -bit binary string  $K_{\text{exp}}$  as an output.

In [27], two approaches are discussed to construct randomness expanders: counter and feedback encryption mode style. HKDF uses a randomness expander that is defined using a key feedback mode over HMAC. HMAC in key feedback mode performs slow due to the frequent rekeying in HMAC calls (consequence of using feedback mode) and double hashing per HMAC call (consequence of using HMAC). In order to avoid these speed breakers, while maintaining 128-bit security level, we turn to the counter mode-like approach and consider some built-in expanding primitives such as a  $\text{PRF}_s$  with  $s \geq 2$  as the underlying primitive. The known counter mode CTR\$ (CTR with random IV) [39] is a good randomness expander, however, in its original form it accommodates  $\gamma$ s only if they are chosen uniformly at random and are of size  $n$  bits (where  $n$  is the underlying block size).

Our proposal FExp is a highly efficient randomness expander that can accommodate arbitrary  $\gamma$ s of size  $2n$  bits. Further, unlike CTR\$ mode which can only provide  $n/2$  bits of security under

the indistinguishability from random strings IND\$ [39] notion, we show that FExp mode (for a given secret random key) provides full  $n$ -bit security under the same security notion. We adjust the notion under the syntax of a randomness expander scheme and denote it by gexp below.

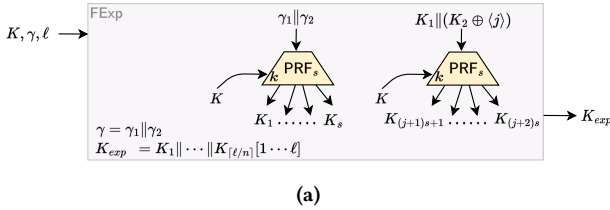
**gexp SECURITY.** The security of a randomness expander or in short, Exp scheme  $\Pi$  is defined with the help of the games **gexp-real** and **gexp-ideal** in Fig. 4. The security of  $\Pi$  is measured as the indistinguishability of its outputs from random strings in a chosen input attack. More precisely, given  $\Pi$  and an adversary  $\mathcal{A}$  who interacts with either **gexp-real** or **gexp-ideal**, we define  $\mathcal{A}$ 's advantage at breaking the gexp security of  $\Pi$  as:

$$\text{Adv}_{\Pi}^{\text{gexp}}(\mathcal{A}) = \left| \Pr[\mathcal{A}^{\text{gexp-real}} \Rightarrow 1] - \Pr[\mathcal{A}^{\text{gexp-ideal}} \Rightarrow 1] \right|.$$

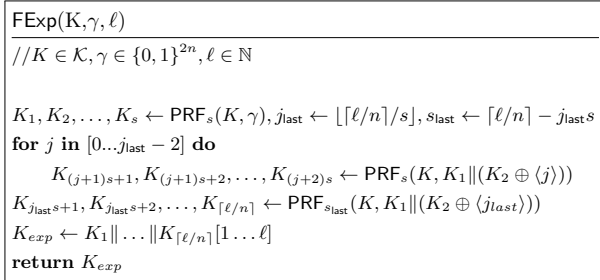
### 6.1 FExp: A PRF<sub>s</sub>-based Randomness Expander

We provide a definition for our randomness expansion scheme based on an expanding (fixed output length) PRF PRF<sub>s</sub>.

**DEFINITION OF FExp.** For a fixed expanding PRF PRF<sub>s</sub> :  $\mathcal{K} \times \{0, 1\}^{2n} \rightarrow \{0, 1\}^{sn}$  with  $s \geq 2$ , FExp takes in a key  $K \in \mathcal{K}$ , the auxiliary info  $\gamma \in \{0, 1\}^{2n}$  and the desired output length  $\ell \in \mathbb{N}$  as inputs. It then uses the PRF<sub>s</sub> as shown in Fig. 3a and outputs the key-stream  $K_{\text{exp}} \in \{0, 1\}^{\ell}$ .



(a)



(b)

**Figure 3: (a) FExp mode of randomness expansion. Here  $K$  denotes the secret key to the PRF<sub>s</sub> and  $\langle j \rangle$  is a suitable  $n$ -bit binary encoding of  $j$  (e.g.,  $0^{n-\lceil \log_2(j) \rceil} \| j$ ). Here each concatenation separates  $n$ -bit (block size) strings. (b) FExp mode pseudocode.**

We give a formal statement of the FExp security and support it with a security proof. The formal security claim is stated in Theorem 6.1.

**THEOREM 6.1 (SECURITY OF FExp).** Let PRF<sub>s</sub> be an expanding pseudorandom function with a secret and uniform random key  $K \in \mathcal{K}$  and  $s \geq 2$ . Then for any adversary  $\mathcal{A}$  who makes at most  $q$  FExp queries such that the total number of PRF<sub>s</sub> calls induced by all the

queries is at most  $\sigma = \sum_{i=1}^q \ell_i$  with  $\ell_i$  being the output length (in  $sn$ -bit blocks) of  $i^{\text{th}}$  query, we have

$$\text{Adv}_{\text{FExp}}^{\text{gexp}}(\mathcal{A}) \leq \text{Adv}_{\text{PRF}_s}^{\text{prf}}(\mathcal{B}) + \frac{2q(\sigma - q)}{2^{2n}}$$

for some adversary  $\mathcal{B}$  who makes at most  $\sigma$  queries, and runs in time given by the running time of  $\mathcal{A}$  plus  $\gamma_0 \cdot \sigma$  for some constant  $\gamma_0$ .

Game <b>gexp-real</b>	Game <b>gexp-ideal</b>
$// K \xleftarrow{\$} \{0, 1\}^k$	$// K \xleftarrow{\$} \{0, 1\}^k$
<b>Oracle</b> $\mathcal{E}(\gamma, \ell)$	<b>Oracle</b> $\mathcal{E}(\gamma, \ell)$
<b>return</b> $\Pi(K, \gamma, \ell)$	<b>return</b> $\text{RF}_{K, \gamma}[1 \dots \ell]$
$b \leftarrow \mathcal{A}^{\mathcal{E}}$	$b \leftarrow \mathcal{A}^{\mathcal{E}}$
<b>return</b> $b$	<b>return</b> $b$

**Figure 4: Games gexp-real and gexp-ideal defining the security of an Exp scheme  $\Pi$ . Here  $\text{RF}_{K, \gamma}$  is a function (independently sampled for every  $(K, \gamma)$ ) that outputs arbitrary many uniform random bits.**

**PROOF OF THEOREM 6.1.** We first replace PRF<sub>s</sub>( $K, \cdot$ ) with a uniformly sampled random function  $f(\cdot) \xleftarrow{\$} \text{Func}(2n, sn)$  and let FExp[ $f$ ] denote the FExp mode that uses  $f$  instead of PRF<sub>s</sub>, which yields

$$\text{Adv}_{\text{FExp}[\text{PRF}_s]}^{\text{gexp}}(\mathcal{A}) \leq \text{Adv}_{\text{PRF}_s}^{\text{prf}}(\mathcal{B}) + \text{Adv}_{\text{FExp}[f]}^{\text{gexp}}(\mathcal{A}).$$

Let us consider that  $\mathcal{A}$  makes at most  $q$  FExp queries with  $i^{\text{th}}$  query containing  $\ell_i$   $f$  calls and hence calling  $f$  for total  $\sigma = \sum_{i=1}^q \ell_i$  many times. Clearly, by definition of  $f$ , we know that all the output bits are random and uniformly distributed as long all the queried  $\sigma$  many inputs to  $f$  are unique. In other words, if all the queried  $f$  inputs in query  $i$  are denoted by the ordered multiset  $Q_i = \{x_j^i\}_{j=1}^{\ell_i} = \{\gamma^i, K_1^i \| K_2^i, K_1^i \| (K_2^i \oplus \langle 1 \rangle), \dots, K_1^i \| (K_2^i \oplus \langle \ell_i - 2 \rangle)\}$  then we have

$$\text{Adv}_{\text{FExp}[f]}^{\text{gexp}}(\mathcal{A}) \leq \Pr[\exists(i, j) < (i', j') \text{ such that } x_j^i = x_{j'}^{i'}]. \quad (3)$$

#### CASE ANALYSIS.

**Case 1** [When  $j = j' = 1$ ]. Under this case, all  $(x_j^i, x_{j'}^{i'})$  pairs are defined as  $(\gamma^i, \gamma^{i'})$  with  $\gamma^i \neq \gamma^{i'} \forall i \neq i'$  and hence  $\Pr[x_j^i = x_{j'}^{i'}] = 0$ .

**Case 2** [When  $j = 1 \vee j' = 1$  but  $j \neq j'$ ]. Under this case, we always have either  $x_j^i = K_1^i \| (K_2^i \oplus \langle j \rangle)$  or  $x_{j'}^{i'} = K_1^{i'} \| (K_2^{i'} \oplus \langle j' \rangle)$  and since each one of  $K_1^i, K_2^i, K_1^{i'}$  and  $K_2^{i'}$  are outputs of a uniform random function ( $f$ ), they are uniformly distributed over  $\{0, 1\}^n$ . Therefore,  $\Pr[x_j^i = x_{j'}^{i'}] = 1/2^{2n}$ . W.l.o.g., let us assume that  $j = 1$  and then there are total  $q$  and at most  $\sigma - q$  many choices to pick  $(i, j)$  and  $(i', j')$ , respectively.

**Case 3** [When  $i = i'$  and  $j \neq 1 \wedge j' \neq 1$ ]. Under this case, all  $(x_j^i, x_{j'}^{i'})$  pairs are defined as  $(K_1^i \| (K_2^i \oplus \langle j \rangle), K_1^i \| (K_2^i \oplus \langle j' \rangle))$  with  $j \neq j'$  hence  $\Pr[x_j^i = x_{j'}^{i'}] = 0$ .



**Case 4** [When  $i \neq i'$  and  $j \neq 1 \wedge j' \neq 1$ ]. Under this case, we always have  $x_j^i = K_1^i \parallel (K_2^i \oplus \langle j \rangle)$  and  $x_{j'}^{i'} = K_1^{i'} \parallel (K_2^{i'} \oplus \langle j' \rangle)$ . Hence, we can write

$$\Pr[x_j^i = x_{j'}^{i'}] = \Pr[(K_1^i \oplus K_1^{i'}) \parallel (K_2^i \oplus K_2^{i'}) = 0^n \parallel (\langle j \rangle \oplus \langle j' \rangle)] \quad (4)$$

Now, since each one of  $K_1^i, K_2^i, K_1^{i'}$  and  $K_2^{i'}$  are outputs of a uniform random function ( $f$ ), they are uniformly distributed over  $\{0, 1\}^n$  and hence  $\Pr[x_j^i = x_{j'}^{i'}] = 1/2^{2n}$ . Additionally, we observe from Eqn. 4 with  $\langle c \rangle = \langle j \rangle \oplus \langle j' \rangle$  that there are total  $\sigma - q$  and at most  $q$  many choices for  $(i, c)$  and  $i'$ , respectively.

Combining all these results into Eqn. 3 gives us

$$\text{Adv}_{\text{FExp}[f]}^{\text{gexp}}(\mathcal{A}) \leq \frac{2q(\sigma - q)}{2^{2n}}$$

and hence the result of Theorem 6.1.  $\square$

## 7 SECURITY ANALYSIS OF Skye

The security of a KDF scheme depends on the properties of the source of initial key material (IKM, Def. 3.1) from which  $Z$  is chosen. We refer the reader to [27] for various examples. Although the definition there does not specify the inputs to the  $\Sigma$  algorithm, it provides a pair  $(Z, C_Z)$  where  $Z$  (the sample set) represents the secret IKM, and  $C_Z$  is a set of some auxiliary knowledge about  $Z$  (or its distribution). This auxiliary information is available to the attacker and can be used in the KDF security analysis. In other words, a KDF is “secure” on inputs  $Z$  even when the value  $C_Z$  is available to the attacker. A Diffie-Hellman value  $Z$  will consist of the value  $g^{xy}$  while  $C_Z$  could represent the set  $\{p, q, g, g^x, g^y\}$ . In a different application, say a random number generator that works by hashing samples of system events in a computer, the value  $C_Z$  may include some of the sampled events used to generate  $Z$ .

Towards defining the security of Skye, we recall the min-entropy and  $m$ -entropy source definitions from [27].

**DEFINITION 7.1 (MIN-ENTROPY [27]).** A probability distribution  $\mathcal{X}$  has min-entropy (at least)  $m$  if for all  $a$  in the support of  $\mathcal{X}$  and for a random variable  $X$  drawn according to  $\mathcal{X}$ ,  $\Pr(X = a) \leq 2^{-m}$ .

**DEFINITION 7.2 ( $m$ -ENTROPY SOURCE [27]).** We say that  $\Sigma$  is a statistical  $m$ -entropy source if for all  $s$  and  $a$  in the support of the distribution  $\Sigma$ , the conditional probability  $\Pr(Z = s | C_Z = a)$  induced by  $\Sigma$  is at most  $2^{-m}$ . We say that  $\Sigma$  is a computational  $m$ -entropy source (or simply an  $m$ -entropy source) if there is a statistical  $m$ -entropy source  $\Sigma'$  that is computationally indistinguishable from  $\Sigma$ .

**CCS SECURITY.** We follow the original CCS definition from [15, 27] for a secure KDF with an  $m$ -entropy source  $\Sigma$  in the adaptive chosen context information ( $\gamma$ ) model with a single salt. The CCS formalization for an  $m$ -entropy source makes use of the **ccs-real** and **ccs-ideal** security games (a simplified version is provided in Fig. 8, App. A). The CCS-security is then defined as the indistinguishability of the KDF generated outputs  $K_{kdf}$  from truly random strings under a chosen input attack. More precisely, given a KDF scheme  $\Pi$  and an adaptive adversary  $\mathcal{A}$ , who interacts with either **ccs-real** or **ccs-ideal**, the  $\mathcal{A}$ 's advantage at breaking the CCS security of  $\Pi$  is defined as  $\text{Adv}_{\Pi}^{\text{CCS}}(\mathcal{A}) = \left| \Pr[\mathcal{A}^{\text{ccs-real}} \Rightarrow 1] - \Pr[\mathcal{A}^{\text{ccs-ideal}} \Rightarrow 1] \right|$ .

We now give the Skye[lsb, PRF $_s$ ] security statement in Theorem 7.1.

**THEOREM 7.1 (SECURITY OF Skye).** Let  $\text{PRF}_s : \mathcal{K} \times \{0, 1\}^{2n} \rightarrow \{0, 1\}^{sn}$  be an expanding pseudorandom function with  $s \geq 2$  and key size  $k \leq 212$  bits. Let  $\Sigma = (Z, C_Z)$  be an input source. Let  $Z$  is sampled from a set containing secret, random, independent and fresh X3DH handshake outputs computed over the group  $G$ , where  $G$  is defined as the cyclic subgroup of Curve25519 using the base point  $x = 9$  [10]. Then, for all adversaries  $\mathcal{A}$  who make  $q$  Skye queries in at most  $\sigma = \sum_{i=1}^q \ell_i$  PRF $_s$  calls with  $\ell_i$  being the output length (in  $sn$ -bit blocks) of  $i^{\text{th}}$  query, we have

$$\text{Adv}_{\text{Skye}[\text{PRF}_s]}^{\text{CCS}}(\mathcal{A}) \leq \text{Adv}_{\text{PRF}_s}^{\text{prf}}(\mathcal{B}) + \text{Adv}_G^{\text{ecddh}}(C) + \frac{q}{2^{128}} + \frac{2q(\sigma - q)}{2^{2n}}$$

for adversaries  $\mathcal{B}$  and  $C$  making at most  $\sigma$  and  $4q$  PRF $_s$  and ECDDH (over  $G$ ) queries, respectively and running in time given by the running time of  $\mathcal{A}$  plus  $\gamma_0 \cdot \sigma$  for some constant  $\gamma_0$ .

**PROOF OF THEOREM 7.1.** Let us first define an event  $E$  which says that the key fed to the underlying FExp component of the Skye construction over a total of  $q$  queries is secret and indistinguishable from a uniform random binary string. Let us now recall from Sec. 6.1 that if the key to the FExp construction is secret and indistinguishable from a uniform random binary string then the outputs of FExp are independent and indistinguishable from uniform random binary strings (of same length) with adversarial advantage as defined in Theorem 6.1. Further, one can also note from the security definition of CCS that the only difference between the real and ideal CCS games w.r.t. Skye is the corresponding outputs being uniform random or not. This as described above under the event  $E$  is upper bounded by  $\text{Adv}_{\text{FExp}[\text{PRF}_s]}^{\text{gexp}}(\mathcal{A}')$  for some adversary  $\mathcal{A}'$  against FExp that uses  $\mathcal{A}$  (restricted under the event  $E$ ) as a subroutine. Now, for Skye with a source group  $G$  defined as in Theorem 7.1 (over Curve25519) we have that under the ECDDH assumption on  $G$ ,  $\Pr(\neg E)$  is upper bounded by  $q \cdot \text{SD}(\text{DExt}_{\text{lsb}}(U_Z), U_k)$  and thus

$$\begin{aligned} \text{Adv}_{\text{Skye}[\text{PRF}_s]}^{\text{CCS}}(\mathcal{A}) &\leq \text{Adv}_G^{\text{ecddh}}(C) + q \cdot \text{SD}(\text{DExt}_{\text{lsb}}(U_Z), U_k) \\ &\quad + \text{Adv}_{\text{FExp}[\text{PRF}_s]}^{\text{gexp}}(\mathcal{A}') \\ &\leq \text{Adv}_{\text{PRF}_s}^{\text{prf}}(\mathcal{B}) + \text{Adv}_G^{\text{ecddh}}(C) + \frac{q}{2^{128}} + \\ &\quad + \frac{2q(\sigma - q)}{2^{2n}}. \end{aligned}$$

The second inequality above is derived from Theorem 6.1 and Corollary 5.2.1 which states that for  $k \leq 212$ ,  $\text{DExt}_{\text{lsb}}$  is a  $(U_Z, 2^{-128})$ -deterministic extractor. This completes the proof of Theorem 7.1.  $\square$

## 8 REAL-WORLD APPLICATIONS OF Skye

From the design description of Skye[ $f$ , PRF $_s$ ], we know that it can be used (with any expanding PRF PRF $_s$ ) as a general KDF for any input source as long as a weak extractor  $f$  is known for that input source.

We use the ButterKnife [7] as the PRF $_s$  in Skye with  $n = 128$  bits and  $s = 8$  due to its (1) large expansion ( $s = 8$ ); (2) well-studied AES internal structure; (3) the support of AES native instructions (NI)

on all supporting processors and (4) efficient key-scheduling [7] for performance gains.

We stress that ButterKnife is not only fast but also secure with 128-bit PRF-security. It comes with dedicated cryptanalysis, a generic proof of security [7], and is backed by further cryptanalytic results for AES-PRF [20, 37] and Deoxys-BC [16, 25, 29, 43].

We also note that the choice of  $\text{PRF}_s$  is independent of the KDF input source or its application. Any other secure PRF can also be used instead of ButterKnife. Like for e.g., HMAC or AES-PRF. However, if the used PRF is *expanding* and *parallelizable* like ButterKnife, we get additional performance advantages (as shown in Sec. 9).

A significant number of KDF applications such as the popular instant messaging apps Signal [34], Signal protocol-based WhatsApp [33], Facebook Messenger [2], Skype [30], Allo [32], Status [5], Secure Chat, Viber and Forsta as well as the Blockchain-based-X3DH [40] for IoTs are based on X3DH or use some multi-DH input sources for key derivation. To illustrate the performance gains of Skye in these multi-DH source applications, in this section we consider Signal as a target application and provide a full performance analysis of Skye as a standalone KDF and when used inside Signal.

We instantiate  $f$  with the least significant bits function  $\text{lsb}$ , as we established its optimality and security for multi-DH sources in Sec. 5.2.

### 8.1 KDF calls in the Signal Protocol

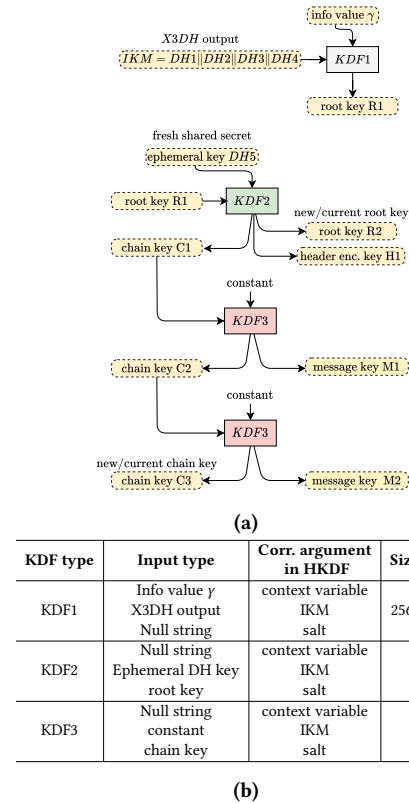
The security in the Signal Protocol is due to the use of the Double Ratchet algorithm [34], the X3DH (triple Elliptic-curve Diffie–Hellman) handshake [35], the key derivation function HKDF [27], and an AEAD mode. Signal Protocol is typically implemented with Curve25519 [10], AES-256, and HMAC-SHA256 as its underlying cryptographic primitives.

Signal as defined in [34]. It makes KDF calls for three different purposes (see Fig. 5a): 1. to generate a **root key** from the X3DH outputs and the info value  $\gamma$ ; 2. to generate a **chain key**, a **header encryption key** and a **new root key** from the present root key and an **ephemeral DH shared secret key**; and 3. to generate a **message key** and a **new chain key** from the old chain key and a predefined constant. The three processes are realized by KDF1, KDF2 and KDF3, respectively, which differ in the types and sizes of inputs and outputs. KDF1 is used first and only once in a session (between two users) to generate the initial root key. When a user sends  $s > 0$  many concurrent messages to another user, the root key is used along with a fresh DH shared secret as an input to KDF2 to produce a chain key and update the root key. Then,  $s$  many iterative KDF3 calls are made with the  $i^{\text{th}}$  chain key as input to generate the  $i^{\text{th}}$  message keys and the  $i + 1^{\text{th}}$  chain key for  $1 \leq i \leq s$ . Fig. 5a illustrates the sequence of calls with its input and output values.

Current implementation libraries [1] of Signal use HKDF [27] for each of these KDF calls. When instantiated with HKDF, the type and size of each input of Signal’s KDF calls in terms of HKDF’s arguments is provided in Table 5b (for 128-bit secure version).

**Skye AS A MORE EFFICIENT ALTERNATIVE TO HKDF FOR SIGNAL.**

**I.** For all KDF1 calls in Signal, the salts for HKDF are set to a constant and the underlying SHA-256 is modeled as a random oracle



**Figure 5: (a) The three types of KDF calls in the Signal Protocol. Here IKM represents the initial key material which is defined as the concatenation of (3 or 4 (depending upon the availability of the corresponding one-time prekeys for DH4, see [34]) DH shared secrets from the X3DH handshake. (b) Types and sizes of Signal’s KDF’s inputs in terms of HKDF-HMAC-SHA-256’s arguments.**

(RO) to achieve the claimed security level (see Lemma 10 and its following paragraphs in [27]). While this is a requirement for general HKDF input sources, it is possible to relax the RO requirement. Further, the reliance on specific randomness sources aids towards a simplified and efficient extraction KDF phase in Signal. In the expansion phase on the other hand, one can benefit both security- and efficiency-wise from an in-built expanding cryptographic primitive with pseudorandom properties, such as the ones captured by the forkcipher [8] (e.g., ForkSkinny [8]; based on the lightweight ISO standard SKINNY [9]), multiforkcipher [6] and expanding PRF [7] (e.g., Butterknife [7]; based on AES and Deoxys-BC [25]) notions. Note that HKDF currently uses a “compressing” primitive instead (HMAC-SHA-256) in a key feedback mode for the expansion phase.

**II.** Fig. 5a shows that only KDF1 requires the use of both the randomness extraction and expansion. KDF2 and KDF3 are made over pseudorandom keys that are generated by the prior KDF calls. Specifically, the root keys that are fed to KDF2 calls are pseudorandom keys as they are generated by the KDF1 calls and the chain keys

that are fed to KDF3 calls are pseudorandom keys as they are generated by the KDF2 calls. In such cases, the extraction phase becomes redundant and the HKDF for KDF2 and KDF3 can be replaced by calling a fast and secure PRF (using the pseudorandom keys) that gives the required expansion.

Addressing these points, we propose Skye[ButterKnife]. We replace HKDF with Skye[ButterKnife] as a KDF for KDF1 and as its expanding component FExp[ButterKnife] for the KDF2 and KDF3 calls. As a consequence, we improve the performance of Signal at full 128-bit security level. For simplicity, we limit our discussion to the 128-bit secure variant. Our results can be analogously lifted to 224-bit security level under suitable instantiation.

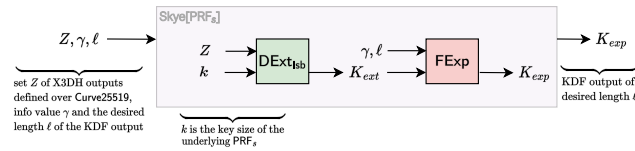
For the sake of completeness, we hereby provide the simplified description of Skye in the context of Signal and Signal-like applications (for general description, see Sec. 4).

**Skye IN SIGNAL.** The inputs to Skye in Signal are:

- the set  $Z$  of  $v = 3$  or  $4$  (depending on the available one time prekeys, see [34]) independent samples  $DH_i$  ( $i = 1, \dots, v$ ) of Diffie-Hellman (DH) shared secrets from the source group  $G$  defined over Curve25519 where each sample is of length  $2n$ ;
- the auxiliary info  $\gamma \in \{0, 1\}^{2n}$ ;
- $\ell \in \mathbb{N}$  denoting the desired output length;

Skye outputs  $K_{exp} \in \{0, 1\}^\ell$  as shown in Fig. 6. DExt<sub>lsb</sub> takes the inputs  $Z$  and a pre-fixed integer  $k$  denoting the output size of DExt<sub>lsb</sub> and produces a  $k$ -bit output  $K_{ext}$ . Then,  $K_{ext}$  together with  $\gamma$  and  $\ell$  are input to the FExp. FExp uses a fixed PRF<sub>s</sub> (in our case ButterKnife) to produce the  $\ell$ -bit  $K_{exp}$ . DExt<sub>lsb</sub> and FExp are given in Fig. 2 and 3a, respectively with pseudocode of FExp provided in Fig. 3b.

Although in Fig. 6  $k$  appears as an external input to DExt<sub>lsb</sub>, for a fixed PRF<sub>s</sub> it is tied to its keysize (in our case  $k = 128$ ). Hence,  $k$  is not a part of the global inputs to Skye[ $f$ , PRF<sub>s</sub>]. The parameter sizes of DExt<sub>lsb</sub> and FExp[ButterKnife] are described in Table 1.



**Figure 6: Skye KDF in the context of Signal. See Fig. 2 and 3a for the internals of functions DExt<sub>lsb</sub> and FExp.**

Parameters	Argument types	Sizes
set $Z$	DExt <sub>lsb</sub> input	$256 \times (3 \text{ or } 4) \text{ bits}$
$K_{ext}$	DExt <sub>lsb</sub> output	$k = 128 \text{ bits}$
info value $\gamma$	FExp input	256 bits
$K_{exp}$	FExp output	$\ell \text{ bits}$

**Table 1: Parameter sizes of input-output arguments in DExt<sub>lsb</sub> and FExp components of Skye[lsb, ButterKnife].**

## 9 SOFTWARE PERFORMANCE OF Skye

In this section, we present the performance evaluation of our implementation of Skye in comparison to HKDF both in an isolated setting and within the current Signal Protocol [1]. We use the Rust implementation of HKDF from the current Signal implementation. Skye’s implementation is also done in Rust, and it is instantiated with ButterKnife. Our measurements are on the x86\_64 platform and we consider performances both with and without using relevant instruction set extensions (AES-NI and SHA-NI). All measurements were performed with AMD Ryzen 7 5800X CPU.

### 9.1 Isolated Performance

To provide a direct comparison, we measured the execution time<sup>1</sup> of sending up to  $n$  messages in a sequence in Signal. This benchmark includes the computation of KDF1, KDF2, and then  $n$  executions of KDF3, but not the time to compute the inputs to the KDF (i.e., they are constants), nor the time to perform any follow-up operations like message encryption with the derived key. On average, our implementation with enabled AES-NI and SHA-NI extensions achieves at least 91% (or  $\geq 11x$ ) speedup, while without the extensions at least 76% (or  $\geq 4x$ ) speedup. On platforms that support AES-NI but not SHA-NI<sup>2</sup>, the speedup is at least 98% (or  $\geq 47x$ ) on average, as HKDF cannot utilize AES-NI, while Skye instantiated with ButterKnife, it can. Full results are shown in Table 2.

$n$	Skye with NI	HKDF with NI	Skye without NI	HKDF without NI
1	157	1739	1794	7404
2	218	2595	2666	11014
3	273	3459	3542	14640
4	336	4318	4414	18275
5	389	5179	5289	21901
6	453	6035	6164	25535
7	505	6894	7037	29160
8	570	7753	7919	32785
9	622	8617	8787	36406
10	688	9470	9667	40036

**Table 2: The mean time (in ns) required to generate  $n$  message keys using Skye and HKDF with and without support for AES-NI and SHA-NI extensions. The measurements were repeated  $10^4$  times.**

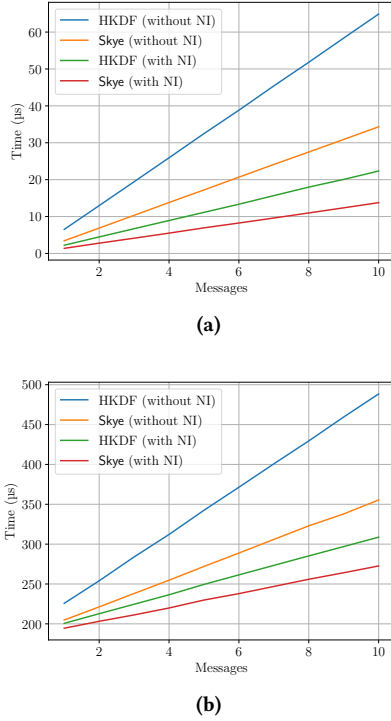
### 9.2 Performance within Signal

We integrated Skye within the current Signal implementation [1]. First, we present the results from the extended benchmarks of the original Signal code. In these settings, we assume that session initialization has already been performed (this includes the one time cost of X3DH computation and KDF1 call), and we measure the subsequent communication.

The Signal library also supports group messaging however, each group message is handled as a direct message to each receiver in

<sup>1</sup>We measure wall-clock time with ns precision to be consistent with benchmarks present in the Signal implementation.

<sup>2</sup>SHA-NI support was released for public markets in 2017-18 with Intel’s Goldmont microarchitecture. All processors and devices before that and many after that do not have the support for SHA-NI.



**Figure 7: (a) The mean time (in  $\mu s$ ) required to encrypt  $n$  messages in the Signal Protocol implementation using Skye and HKDF. (b) The mean time (in  $\mu s$ ) required to encrypt (on sender side) then decrypt (on receiver side)  $n$  messages, and then encrypt (on receiver side) and decrypt (on sender side) another  $n$  messages in the Signal Protocol implementation using Skye and HKDF. In both cases, the measurements were performed for variants with and without AES-NI and SHA-NI support. The encrypted messages were of less than 16 characters. All measurements were repeated  $10^3$  times.**

the group, i.e., if there are  $N$  members in the group, Signal client sends  $N$  messages individually encrypted with the derived ratchet key of each participant. Hence, we focus on two-party messaging for analyzing performance in both direct and group messaging.

**UNIDIRECTIONAL MESSAGING:** Here one party sends messages to another in a sequence without any reply in between. Fig. 7a presents the time it takes for the party to encrypt  $n$  messages.<sup>3</sup> The relative speedup is independent of  $n$  and on average equal to 38% with NI instructions, 47% without NI instructions, and 64% with AES-NI but not SHA-NI<sup>4</sup>. For extended experiment results, see Table 4 in App. E.

**BIDIRECTIONAL MESSAGING:** The previous setting does not include the cost of DH-ratchet, hence we consider: a party sends  $n$  messages (in sequence) to another party, which then replies with

<sup>3</sup>We use one block test messages of size 128 bits (i.e. 16 characters).

<sup>4</sup>Note that in this case AES-NI and SHA-NI were also used for message encryption and authentication so this number (64%) cannot be computed from the figure and is based on the full measurement.

$n$  messages. In Fig. 7b, we present the time required to encrypt and decrypt the sent and received messages, respectively. Since the DH computation cost diminishes the impact of the rest of the computation (especially when the KDF calls are a few), we gain only around 3% with single message i.e.  $n = 1$  with NI instructions, but (as expected) with more messages, the speedup improves and comes closer to the unidirectional case. For example, the speedup quickly improves here to 12% for  $n = 10$ . Similarly, without NI, we gain a speedup of around 9% with  $n = 1$ , and with  $n = 10$ , up to 27%. Lastly, with the partial extension support (with AES-NI but not SHA-NI) we get 11% with  $n = 1$  that goes to 36% with  $n = 10$ . For extended experiment results, see Table 5 in App. E.

**SPEEDUP DUE TO THE EXPANDING PRF.** In the unidirectional experiment, 81%-93% of the total performance gain is due to the use of an expanding PRF (ButterKnife) in place of a compression function and the rest 19%-7% is due to the replacement of KDF2 and KDF3 calls by FExp. In the bidirectional experiment, with 10 messages, 78%-86% of the total performance gain is due to the use of an expanding PRF and the rest 22%-14% is due to the KDF2 and KDF3 replacements. These results illustrate the explicit performance advantage of adding an expanding PRF.

## 10 DISCUSSION AND FUTURE RESEARCH

**Curve448 AND 224-BIT SECURITY.** Signal with Skye targets 128-bit security with Curve25519. The 224-bit secure version of Signal is based on Curve448. This puts forward the question of whether the Skye syntax and similar security proof arguments apply for the 224-bit version. The general study of the underlying components of Skye is directly applicable for constructing 224-bit secure Skye. A 224-bit secure Skye can be defined with the same syntax, and similar security arguments when Curve25519 is replaced by Curve448. However, we avoid a separate security treatment of Skye with Curve448 given that no PRF with key security  $> 128$  bits, input size  $> 256$  bits and no idealized assumptions exists. We leave the research of finding efficient expanding PRFs with at least 224 bits security in the standard model as an open problem.

**NETWORK LATENCY AND IOT DEVICES.** The changes we propose in this work improve performance of the Signal cryptographic component, but the impact on the overall performance of modern smartphones and laptops might reduce when network latency is included. Nonetheless, this tradeoff differs for all various of devices and network setups. Low-performance hardware, e.g., IoT devices or wireless sensor networks, could greatly benefit from our improvements, enabling a wider deployment of the Signal Protocol.

**ENERGY EFFICIENCY.** Although we present improved time efficiency, an associated benefit is the lower energy consumption. Energy savings are relevant even in settings where the network latency overshadows the performance.

**SECURITY UNDER COMPROMISED X3DH SAMPLES.** To avoid the RO-model, our present Skye analysis in Signal assumes at least 3 uncompromised X3DH keys. That results in stronger standard model security when compared to HKDF. In the case when there are at most 2 uncompromised X3DH keys (and assuming the user doesn't update the compromised ones), security can be maintained

albeit by modeling the  $\text{PRF}_s$  as a RO (similarly to the treatment in [19]), a subject to separate analysis.

**Skye APPLICATIONS BEYOND SIGNAL AND DH SOURCES.** Skye is based on the generic  $\text{DExt}_f$ , hence its security is not limited to Curve25519 points or DH sources. Realizing the potential of Skye beyond DH sources, and finding corresponding weak extractors  $f_s$ , is also an important direction for future works. This entails exploring applications, such as the present Message Layer Security (MLS) [3] and Post-Quantum Signal projects such as [12] where KDF calls are made to extract uniformly random keys from shared secrets. Further performance evaluation of Skye for applications beyond Signal is also an interesting future work.

**ON INSTANTIATION OF  $\text{PRF}_s$ .** In this work we use the concrete instantiation of ButterKnife. Yet, our design can accommodate any  $\text{PRF}_s$ . For example, ForkSkinny is also a suitable candidate but then the gains, when compared with ButterKnife, would be lessened. The PRF security reduces to half, 64-bits, and the performance as well, due to the fact that ForkSkinny performs slower than ButterKnife and does not have NI support on regular platforms.

**DExt BEYOND LINEAR SYSTEMS.** The number of extracted bits from DExt and the security might be improved for the same inputs when non-linear multivariate equations (with the extra cost of field multiplications) are used. Studying such systems is out of the scope and motivation of this paper and we leave the design and analysis of further generalized and efficient DExt-like extensions as an open problem.

## ACKNOWLEDGEMENTS

This work was supported by CyberSecurity Research Flanders with reference number VR20192203. Amit Singh Bhati and Bart Preneel were supported in part by the Research Council KU Leuven C1 on Security and Privacy for Cyber-Physical Systems and the Internet of Things with contract number C16/15/058 and by the Flemish Government through FWO Project G.0835.16 A security Architecture for IoT.

## REFERENCES

- [1] [n.d.]. Signal Protocol software libraries. ([n.d.]). Github [accessed on 05/02/2023], <https://github.com/signalapp/>.
- [2] 2016. Messenger secret conversations: Technical whitepaper. (2016). [https://fbnewsroomus.files.wordpress.com/2016/07/secret\\_conversations\\_whitepaper-1.pdf](https://fbnewsroomus.files.wordpress.com/2016/07/secret_conversations_whitepaper-1.pdf).
- [3] 2022. The Messaging Layer Security (MLS) Protocol. (2022). Work in Progress, Internet-Draft, draft-ietf-mls-protocol-17, 19 December 2022, <https://datatracker.ietf.org/doc/html/draft-ietf-mls-protocol-17>.
- [4] Joël Alwen, Sandro Coretti, and Yevgeniy Dodis. 2019. The Double Ratchet: Security Notions, Proofs, and Modularization for the Signal Protocol. In *Advances in Cryptology – EUROCRYPT 2019*, Yuval Ishai and Vincent Rijmen (Eds.). Springer International Publishing, Cham, 129–158.
- [5] Corey Petty Oskar Thorén Dean Eigenmann Andrea Piana, Pedro Pombeiro. [n.d.]. Specifications for Status clients - 5/SECURE-TRANSPORT. ([n.d.]). <https://specs.status.im/spec/5>.
- [6] Elena Andreeva, Amit Singh Bhati, Bart Preneel, and Damian Vizár. 2021. 1, 2, 3, Fork: Counter Mode Variants based on a Generalized Forkcipher. *IACR Trans. Symmetric Cryptol.* 2021, 3 (2021), 1–35.
- [7] Elena Andreeva, Benoit Cogliati, Virginie Lallemand, Marine Minier, Antoon Purnal, and Arnab Roy. 2022. Masked Iterate-Fork-Iterate: A new Design Paradigm for Tweakable Expanding Pseudorandom Function. *Cryptology ePrint Archive* (2022).
- [8] Elena Andreeva, Virginie Lallemand, Antoon Purnal, Reza Reyhanitabar, Arnab Roy, and Damian Vizár. 2019. Forkcipher: a New Primitive for Authenticated Encryption of Very Short Messages. In *International Conference on the Theory and Application of Cryptology and Information Security*. Springer, 153–182.
- [9] Christof Beierle, Jérémy Jean, Stefan Kölbl, Gregor Leander, Amir Moradi, Thomas Peyrin, Yu Sasaki, Pascal Sasdrich, and Siang Meng Sim. 2016. The SKINNY family of block ciphers and its low-latency variant MANTIS. In *Annual International Cryptology Conference (CRYPTO)*. Springer, 123–153.
- [10] Daniel J Bernstein. 2006. Curve25519: new Diffie-Hellman speed records. In *International Workshop on Public Key Cryptography*. Springer, 207–228.
- [11] Dan Boneh. 1998. The Decision Diffie-Hellman problem. In *Algorithmic Number Theory*, Joe P. Buhler (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 48–63.
- [12] Jacqueline Brendel, Rune Fiedler, Felix Günther, Christian Janson, and Douglas Stebila. 2022. Post-quantum asynchronous deniable key exchange and the signal handshake. In *Public-Key Cryptography–PKC 2022: 25th IACR International Conference on Practice and Theory of Public-Key Cryptography, Virtual Event, March 8–11, 2022, Proceedings, Part II*. Springer, 3–34.
- [13] Lily Chen, Dustin Moody, Andrew Regenscheid, and Karen Randall. 2019. *SP800-186, Recommendations for discrete logarithm-based cryptography: Elliptic curve domain parameters*. Technical Report. National Institute of Standards and Technology.
- [14] Céline Chevalier, Pierre-Alain Fouque, David Pointcheval, and Sébastien Zimmer. 2009. Optimal Randomness Extraction from a Diffie-Hellman Element. *EUROCRYPT 2009* (2009), 572.
- [15] Chai Wen Chuah, Edward Dawson, and Leonie Simpson. 2013. Key derivation function: the SCKDF scheme. In *IFIP International Information Security Conference*. Springer, 125–138.
- [16] Carlos Cid, Tao Huang, Thomas Peyrin, Yu Sasaki, and Ling Song. 2017. A security analysis of Deoxys and its internal tweakable block ciphers. *IACR Transactions on Symmetric Cryptology* (2017), 73–107.
- [17] Abdoul Aziz Ciss. 2014. Two-sources randomness extractors for elliptic curves. *arXiv preprint arXiv:1404.2226* (2014).
- [18] Abdoul Aziz Ciss and Djiby Sow. 2017. Two-Source Randomness Extractors for Elliptic Curves for Authenticated Key Exchange. In *International Conference on Codes, Cryptology, and Information Security*. Springer, 85–95.
- [19] Katriel Cohn-Gordon, Cas Cremers, Benjamin Dowling, Luke Garratt, and Douglas Stebila. 2017. A Formal Security Analysis of the Signal Messaging Protocol. In *2017 IEEE European Symposium on Security and Privacy (EuroS P)*. 451–466. <https://doi.org/10.1109/EuroSP.2017.27>
- [20] Patrick Derbez, Tetsu Iwata, Ling Sun, Siwei Sun, Yosuke Todo, Haoyang Wang, and Meiqin Wang. 2018. Cryptanalysis of AES-PRF and its dual. *IACR Transactions on Symmetric Cryptology* 2018, 2 (2018).
- [21] Yevgeniy Dodis, Ariel Elbaz, Roberto Oliveira, and Ran Raz. 2004. Improved randomness extraction from two independent sources. In *Approximation, randomization, and combinatorial optimization. Algorithms and techniques*. Springer, 334–344.
- [22] Pierre-Alain Fouque, David Pointcheval, Jacques Stern, and Sébastien Zimmer. 2006. Hardness of Distinguishing the MSB or LSB of Secret Keys in Diffie-Hellman Schemes. In *Automata, Languages and Programming, 33rd International Colloquium, ICALP, Venice, Italy, 2006, Part II (Lecture Notes in Computer Science, Vol. 4052)*, Michele Bugliesi, Bart Preneel, Vladimiro Sassone, and Ingo Wegener (Eds.). Springer, 240–251. [https://doi.org/10.1007/11787006\\_21](https://doi.org/10.1007/11787006_21)
- [23] Edgar N Gilbert. 1952. A comparison of signalling alphabets. *The Bell system technical journal* 31, 3 (1952), 504–522.
- [24] Markus Grassl. 2007. Bounds on the minimum distance of linear codes and quantum codes. Online available at <http://www.codetables.de>. Accessed on 2021-06-25.
- [25] Jérémy Jean, Ivica Nikolić, Thomas Peyrin, and Yannick Seurin. October 2016. Submission to CAESAR : Deoxys v1.41. <http://competitions.cr.yt.to/round3/deoxysv141.pdf>.
- [26] DT Kolyang, Djiby Sow, Abdoul A Ciss, and Hortense Boudjou Tchagnouo. 2017. Two-sources randomness extractors in finite fields and in elliptic curves. *African Journal of Research in Computer Science and Applied Mathematics* 24 (2017).
- [27] Hugo Krawczyk. 2010. Cryptographic extraction and key derivation: The HKDF scheme. In *Annual Cryptology Conference*. Springer, 631–648.
- [28] Hugo Krawczyk, Mihir Bellare, and Ran Canetti. 1997. HMAC: Keyed-hashing for message authentication.
- [29] Ya Liu, Bing Shi, Dawu Gu, Fengyu Zhao, Wei Li, and Zhiqiang Liu. 2020. Improved meet-in-the-middle attacks on reduced-round Deoxys-BC-256. *Comput. J.* 63, 12 (2020), 1859–1870.
- [30] J Lund. [n.d.]. Signal partners with Microsoft to bring end-to-end encryption to Skype. ([n.d.]). <https://signal.org/blog/skype-partnership/>.
- [31] Neil Madden. [n.d.]. What’s the Curve25519 clamping all about? ([n.d.]). <https://neilmadden.blog/2020/05/28/whats-the-curve25519-clamping-all-about/>.
- [32] Moxie Marlinspike. 2016. Open whisper systems partners with Google on end-to-end encryption for Allo. (2016). <https://signal.org/blog/allo/>.
- [33] Moxie Marlinspike. 2016. WhatsApp’s Signal Protocol integration is now complete. (2016). <https://signal.org/blog/whatsapp-complete/>.



- [34] Moxie Marlinspike and Trevor Perrin. [n.d.]. The double Ratchet algorithm, November 2016. ([n.d.]). <https://whispersystems.org/docs/specifications/doublerratchet/doublerratchet.pdf>.
- [35] Moxie Marlinspike and Trevor Perrin. 2016. The X3DH key agreement protocol. *Open Whisper Systems* (2016). <https://signal.org/docs/specifications/x3dh/x3dh.pdf>.
- [36] Ueli Maurer, Krzysztof Pietrzak, and Renato Renner. 2007. Indistinguishability amplification. In *Advances in Cryptology-CRYPTO 2007: 27th Annual International Cryptology Conference, Santa Barbara, CA, USA, August 19-23, 2007. Proceedings* 27. Springer, 130–149.
- [37] Bart Mennink and Samuel Neves. 2017. Optimal PRFs from blockcipher designs. *IACR Transactions on Symmetric Cryptology* (2017), 228–252.
- [38] Trevor Perrin. 2016. The Noise protocol framework. (2016). [noiseprotocol.org](https://noiseprotocol.org).
- [39] Phillip Rogaway. 2011. Evaluation of some blockcipher modes of operation. *Cryptography Research and Evaluation Committees (CRYPTREC) for the Government of Japan* (2011). [https://crossbowertb.github.io/docs/crypto/rogaway\\_modes.pdf](https://crossbowertb.github.io/docs/crypto/rogaway_modes.pdf).
- [40] Armando Ruggeri, Antonio Celesti, Maria Fazio, Antonino Galletta, and Massimo Villari. 2020. BCB-X3DH: a Blockchain Based Improved Version of the Extended Triple Diffie-Hellman Protocol. In *2020 Second IEEE International Conference on Trust, Privacy and Security in Intelligent Systems and Applications (TPS-ISA)*, 73–78. <https://doi.org/10.1109/TPS-ISA50397.2020.00020>
- [41] Hortense Boudjou Tchagnouo and Abdoul Aziz Ciss. 2015. Multi-sources Randomness Extraction over Finite Fields and Elliptic Curve. *arXiv preprint arXiv:1502.00433* (2015).
- [42] Rom Rubenovich Varshamov. 1957. Estimate of the number of signals in error correcting codes. *Doklady Akad. Nauk, SSSR* 117 (1957), 739–741.
- [43] Boxin Zhao, Xiaoyang Dong, and Keting Jia. 2019. New related-tweakey boomerang and rectangle attacks on Deoxys-BC including BDT effect. *IACR Transactions on Symmetric Cryptology* (2019), 121–151.

## A CCS SECURITY GAMES

Game <b>ccs-real</b>	Game <b>ccs-ideal</b>
$//I_\gamma = \{\}$	$//I_\gamma = \{\}$
$//\gamma_c = \perp$ until provided by $\mathcal{A}$	$//\gamma_c = \perp$ until provided by $\mathcal{A}$
$(Z, C_Z) \leftarrow \Sigma // m\text{-entropy source}$	$(Z, C_Z) \leftarrow \Sigma // m\text{-entropy source}$
<b>Oracle</b> $\text{Derive}(\gamma, \ell)$	<b>Oracle</b> $\text{Derive}(\gamma, \ell)$
<b>define</b> $I_\gamma = I_\gamma \cup \{\gamma\}$	<b>define</b> $I_\gamma = I_\gamma \cup \{\gamma\}$
<b>if</b> $\gamma \neq \gamma_c$	<b>if</b> $\gamma \neq \gamma_c$
<b>return</b> $(\Pi(Z, \gamma, \ell), \text{salt})$	<b>return</b> $(\Pi(Z, \gamma, \ell), \text{salt})$
<b>else</b> $//\text{challenge query}$	<b>else</b> $//\text{challenge query}$
<b>return</b> $(\Pi(Z, \gamma, \ell), \text{salt})$	<b>return</b> $(\text{RF}_{Z, \gamma}[1 \dots \ell], \text{salt})$
$(\gamma_c, \ell_c), b \leftarrow \mathcal{A}^{\text{Derive}}$	$(\gamma_c, \ell_c), b \leftarrow \mathcal{A}^{\text{Derive}}$
$//\text{challenge query and the}$	$//\text{challenge query and the}$
$//\text{final output bit of } \mathcal{A}$	$//\text{final output bit of } \mathcal{A}$
<b>if</b> $\gamma_c \in I_\gamma$	<b>if</b> $\gamma_c \in I_\gamma$
<b>return</b> $\perp$	<b>return</b> $\perp$
<b>else</b>	<b>else</b>
<b>return</b> $b$	<b>return</b> $b$

**Figure 8: Games **ccs-real** and **ccs-ideal** defining the CCS-security of a KDF scheme. Here  $\text{RF}_{Z, \gamma}$  is a function (independently sampled for every  $(Z, \gamma)$ ) that outputs arbitrary many uniform random bits and  $\text{salt}$  denotes the internally sampled salt value which is used in the queries  $\text{Derive}(\cdot, \cdot)$  (Note that this salt here can be a uniform random string or fixed to some constant or null string. However, we let its sampling remain general in the security definition).**

## B (SUB-)OPTIMALITY ANALYSIS OF $\text{DExt}_f$

In this section we provide the security analysis of our proposed extractor and its instantiation. We start with the basic definitions that are required for the analyses presented in this section.

**DEFINITION B.1 (CONSISTENT SYSTEM).** *A system of linear equations over a field  $\mathbb{F}$  is called consistent if it has at least one solution in  $\mathbb{F}$ .*

**DEFINITION B.2 (COEFFICIENT MATRIX).** *For a given set  $S$  of linear equations, the coefficient matrix  $C_S$  of  $S$  is defined by a matrix with  $i^{\text{th}}$  row as the coefficients of the variables in the  $i^{\text{th}}$  linear equation of the set  $S$ .*

**DEFINITION B.3 (BINARY LINEAR CODE).** *A binary linear code of length  $n$  and rank  $k$  is a linear subspace  $C$  with dimension  $k$  of the vector space  $\mathbb{F}_2^k$  where  $\mathbb{F}_2$  is the binary field. The vectors in  $C$  are called codewords.*

**DEFINITION B.4 (HAMMING DISTANCE).** *For any binary linear code  $C$ , the Hamming distance between any two codewords in  $C$  is defined as the number of elements in which the codewords differ.*

Any binary linear code  $C$  can be represented by  $[n, k, d]_2$  where  $n$  is the length of a codeword in  $C$ ,  $d$  is the minimum of Hamming distances of all pairs of codewords in  $C$  and  $k$  is the rank of  $C$ . The maximum number of codewords in any  $[n, k, d]_2$  code is denoted by  $A_2(n, d)$ . This implies that  $k \leq \log_2(A_2(n, d))$ .

We now introduce two new definitions called degree of involvement and min-degree of involvement for the upcoming analysis.

**DEFINITION B.5 (DEGREE OF INVOLVEMENT).** *Let  $S = \{E_i | 1 \leq i \leq a\}$  be a consistent system of equations of  $v$  variables over a field  $\mathbb{F}$  and let  $x_0$  be one of these  $v$  variables then the degree of involvement  $\text{dis}_S$  of  $x_0$  is defined as the minimum number of other variables on which the value of  $x_0$  depends in  $S$ .*

To exemplify, consider the following system of one equation with 3 variables over  $\mathbb{R}$ ,  $S_1 = \{x + y + z = 3\}$  where  $V = \{x, y, z\}$ . Clearly, any variable in  $V$  has 2 degrees of involvement as its value depends on the other two variables. Now, as another example, consider the following system of three equations over  $\mathbb{R}$ ,  $S_2 = \{x + y + z = 3, x + y = 2, z = 1\}$ . Unlike to the previous example, in  $S_2$  where  $z$  is now fixed by another equation, we have zero degree of involvement left in the system for  $z$ . Similarly, for  $x$  and  $y$ , the degree of involvement is 1 as they depend on each other.

**DEFINITION B.6 (MIN-DEGREE OF INVOLVEMENT).** *Let  $S = \{E_i | 1 \leq i \leq a\}$  be a consistent system of equations of  $v$  variables defined by the set  $V = \{V_1, V_2, \dots, V_v\}$  over a field  $\mathbb{F}$  then the min-degree of involvement (MDI)  $\text{mdi}$  of  $S$  is defined as the minimum of  $\text{dis}_S(V_i)$  for  $1 \leq i \leq a$ .*

We note that the MDI of a system is not equal to its degree of freedom. In fact, one can show that the degree of freedom for a system  $S$  is equal to the “maximum” of  $\text{dis}_S(x)$  over all variables  $x$  in  $S$ . However, for this work the important extremum on these degrees is the minimum defined above as the MDI of a system. The definition of degree of freedom is provided below for completeness.

**DEFINITION B.7 (DEGREE OF FREEDOM).** *Let  $S = \{E_i | 1 \leq i \leq a\}$  be a system of linear equations of  $v$  variables defined by the set*



$V = \{V_1, V_2, \dots, V_v\}$  over a field  $\mathbb{F}$  and let  $V_f \subset V$  be the largest subset of  $V$  such that for all possible values of the variables of  $V_f$  in  $\mathbb{F}$ , all equations of  $S$  holds. We use the term “free variable” to denote a variable in  $V_f$  and the degree of freedom  $\text{df}$  of  $S$  is defined as  $\text{df}(S) = |V_f|$ .

Next, we state necessary theorems towards the construction of  $\text{DExt}_f$ . We emphasize that to extract optimal randomness from a given input  $a = \{a_i\}_{i=1}^v$  and  $c$ , one would need to find one of the largest system  $S$  of linear equations containing  $v$  variables (each variable corresponding to an element in the set  $a$ ) with binary coefficients and with  $\text{mdi}(S) = c - 1$ . Below in Theorem B.1 we show that this problem is equivalent to finding a  $[v, \log_2(A_2(v, c)), c]_2$  optimal binary linear code which has been considered hard for general values of  $v, c$  and that there is no polynomial-time algorithm that can find a  $[v, \log_2(A_2(v, c)), c]_2$  code for arbitrary values of  $v, c$ .

**THEOREM B.1.** *Let  $\mathcal{S}_c$  be the collection of all consistent systems  $S^j = \{E_i^j | 1 \leq i \leq a_j\}$  of linear equations in  $v$  variables defined by the set  $V = \{V_1, V_2, \dots, V_v\}$  over the binary field  $\mathbb{F}_{2^w}$  (for some integer  $w \geq 0$ ) with binary coefficients and for each  $S \in \mathcal{S}_c$  we have  $\text{mdi}(S) = c - 1$  for some positive integer  $c \leq v$  then for a system  $S^* \in \mathcal{S}_c$  such that  $|S^*| = \max_j \{a_j | S_j \in \mathcal{S}_c\}$  we have*

$$|S^*| = \log_2(A_2(v, c)).$$

Proof of Theorem B.1 is straightforward from the fact that the coefficient matrix  $C_S$  for any system  $S \in \mathcal{S}_c$  can be equivalently seen as a basis set of codewords a.k.a. the generator matrix for a  $[v, |S|, c]_2$  code. Hence, for  $S^*$  as defined, we have  $2^{|S^*|} = A_2(v, c)$  and thus the result of the theorem.

We further note that there exists a good lower bound on  $A_2(v, c)$  called the Gilbert–Varshamov bound [23, 42] (which states that  $\log_2(A_2(v, c)) \geq \lfloor v - \log_2 \sum_{i=0}^{c-2} \binom{v}{i} \rfloor$ ) and proves the existence of a  $[v, |S|, c]_2$  code with  $|S| \geq \lfloor v - \log_2 \sum_{i=0}^{c-2} \binom{v}{i} \rfloor$ , however, there does not exist any deterministic method that can construct a linear code satisfying the Gilbert–Varshamov (GV) bound.

For our work, we settle with a comparatively loose lower bound of  $|S| = b = \lfloor (v - c)/\lceil c/2 \rceil \rfloor + 1$  (which is very close to GV or even optimal for settings where  $v$  is very small and quite loose, otherwise) but with a deterministic algorithm  $\mathcal{E}$  that can construct linear codes satisfying this bound (as shown below in Def. B.8 and Theorem B.2).

**DEFINITION B.8.**  $\mathcal{E}$  is a deterministic algorithm that, when provided with two integers  $v$  and  $c \leq v$ , first computes a coefficient matrix  $C_S$  with its  $ij^{\text{th}}$  entry  $C_S[ij]$  defined as

$$C_S[ij] = \begin{cases} 1 & \text{if } \lceil c/2 \rceil(i - 1) + 1 \leq j \leq \lceil c/2 \rceil(i - 1) + c \\ 0 & \text{o/w} \end{cases}$$

and then for a variable set  $\{V_1, V_2, \dots, V_v\}$  with independent variables over the binary field  $\mathbb{F}_{2^w}$  (for some integer  $w \geq 0$ ), it returns the corresponding consistent system  $S$  of  $C_S$ .

**THEOREM B.2.** *Every system  $S$  returned by the algorithm  $\mathcal{E}$  as defined above in Def. B.8 has  $\text{mdi}(S) = c - 1$ .*

We defer the proof of Theorem B.2 to App. B.1 and provide Table 3 in App. C showing the differences between  $\log_2(A_2(v, c))$

and  $b = \lfloor (v - c)/\lceil c/2 \rceil \rfloor + 1$  for various values of  $v$  and  $c$  up to sample sizes  $v \leq 10$ . From the table, one can infer that our definition of  $b$  can be considered very good for small values of  $v$  and even optimal for  $v \leq 5$  (which is pretty sufficient in practice).

Further, we note that the construction of  $\text{DExt}_f$  allows variable sample sizes, security parameters and output sizes, therefore, for a deterministic and faster execution of this algorithm, we fix  $b$  to  $\lfloor (v - c)/\lceil c/2 \rceil \rfloor + 1$ . However, we recommend that for applications where sample size, security parameter and output size is defined only once, one may use a better or even optimal code (i.e.  $\log_2(A_2(v, c))$  instead of  $b$ ), if exists to extract more randomness with almost the same security. We refer the reader to codetables.de[24] for existing tables of optimal  $[v, \log_2(A_2(v, c)), c]_2$  codes for certain  $v, c$  settings.

In this paper,  $b$  is treated as a positive integer, defined as  $\lfloor (v - c)/\lceil c/2 \rceil \rfloor + 1$  for another positive integers  $v, c$ . Therefore, all results that are defined in terms of  $b$  and  $c$  in this paper are directly applicable to any application for which a different and larger value for  $b$  exists.

For simplicity, we have provided (and used in the rest of this paper) an equivalent definition of the algorithm  $\mathcal{E}$  as a function called  $\text{invXOR}$  (see Def. 5.3).

Equivalence between the two definitions can be easily understood from the fact that for a given variable set  $\{V_1, V_2, \dots, V_v\}$  and their corresponding values as a set  $\{a_1, a_2, \dots, a_v\}$ , each concatenation in the output of  $\text{invXOR}$  function corresponds to an equation in the generated system  $S$  of  $\mathcal{E}$  and vice versa.

## B.1 Proof of Theorem B.2

**PROOF.** Let  $S$  denote a system returned by the algorithm  $\mathcal{E}$  as defined in Def. B.8 and let  $C_S$  denote the corresponding coefficient matrix of  $S$  then we have that the  $ij^{\text{th}}$  entry of  $C_S$  can be defined as

$$C_S[ij] = \begin{cases} 1 & \text{if } \lceil c/2 \rceil(i - 1) + 1 \leq j \leq \lceil c/2 \rceil(i - 1) + c \\ 0 & \text{o/w} \end{cases}$$

Let  $U = \{E_{i_1}, E_{i_2}, \dots, E_{i_x}\}$  denote an arbitrary subset of the system  $S$  with size  $x > 1$ . Clearly, if we show that the combined XOR of all equations in  $U$  always contains at least  $c$  many 1s then we can say that no linear combination of equations in  $S$  can have degree of involvement  $< c - 1$  and hence the claim of the Theorem.

Now, to prove the above statement, we use the following simple approach. Let us first define an indexed set  $U'$  as the sorted version of  $U$  where the equations are sorted by the value of their corresponding first column indices  $j$ s in  $C_S$  with  $C[ij]$  entry as 1. In other words, the sorted set  $U'$  will have the entries in the same order as they are defined in  $S$ . Clearly, the combined XOR of  $U$  will be same as of  $U'$ . Now, one can note that in this definition of  $C_S$ , every row contains at least  $\lceil c/2 \rceil$  many unique 1s entries than others. Hence, the combined XOR of  $U'$  will always have the unique 1s entries of the first and the last equation of  $U'$  which in total will be  $2\lceil c/2 \rceil \geq c$  many 1s.  $\square$

## C A CODE DIFFERENCE TABLE

$c \backslash v$	1	2	3	4	5	6	7	8	9	10
1	0	0	0	0	0	0	0	0	0	0
2		0	0	0	0	0	0	0	0	0
3			0	0	0	1	1	1	1	2
4				0	0	0	1	1	1	1
5					0	0	0	0	0	1
6						0	0	0	0	0
7							0	0	0	0
8								0	0	0
9									0	0
10										0

**Table 3: A table showing the differences between the actual value of  $\log_2(A_2(v, c))$  and  $b$  for various values of  $v$  and  $c$  upto sample sizes  $v \leq 10$ . The optimal values of  $\log_2(A_2(v, c))$  are taken from codetables.de[24].**

## D PRF-PRNG VS CCS SECURITY

Alwen et al. [4] proposed a syntax called PRF-PRNG for KDFs in Signal. In simpler words, a PRF-PRNG takes three inputs - the current state/key  $\sigma$ , an input  $I$  (which in Signal's context will be  $Z||\gamma$ ) and the output length  $\ell$  (left implicit in the original description of [4]) and returns a string  $R$  of length  $\ell$  bits and a new state/key to be updated  $\sigma'$ . Any PRF-PRNG = (P-init, P-up) is consist of two algorithms - P-init which is used to initiate the PRF-PRNG by generating the first state/key  $\sigma$  (for this it may use some preshared secret key) and P-up which is used over the PRF-PRNG inputs as defined above to generate the outputs.

The paper also defines a security notion P for PRF-PRNG schemes which we will refer here as prf-prng for ease of understanding. For Signal's application, the same paper mentions two KDF examples with prf-prng security, namely, HKDF and PRP-then-PRG [4]. We note here that HKDF is prf-prng-secure when its underlying compression function is assumed as RO whereas the proposed PRP-then-PRG is prf-prng-secure when the PRP is initialized with a preshared "uniform" random secret key. This assumption of preshared symmetric key can not be met in practice for two arbitrary parties that are communicating for the first time which makes PRP-then-PRG inapplicable to Signal. We also note that this incompatibility can be countered by assuming the PRP to be an ideal cipher (which is analogous to the RO assumption) and hence initializing the key with some public constant.

**HOW TO DEFINE A PRF-PRNG USING KDFs WITH STANDARD SYNTAX.** We first formally define a simple way of converting a KDF with the standard syntax (as defined in [27]) into a PRF-PRNG and then state the formal claim on relation among the two mentioned KDF security notions in Theorem D.1.

Let  $\Pi(Z, \gamma, \ell)$  and  $\Pi'(Z', \gamma', \ell')$  be two standard syntax KDF functions then we can define a PRF-PRNG as  $\Pi_{pp} = (\text{P-init}, \text{P-up})$  where P-init is just a constant function that returns  $\sigma_0 = 0^s$  for  $s$  being the state size of the PRF-PRNG and for the  $i(\geq 1)^{th}$  query to the

PRF-PRNG, we have

$$\begin{aligned} \text{P-up}(\sigma_{i-1}, I_i = Z_i || \gamma_i, \ell_i) \\ = \begin{cases} \sigma_1 || R_1 = \Pi(Z_1, \gamma_1, \ell_1 + s) & \text{when } i = 1 \\ \sigma_i || R_i = \Pi'(\sigma_{i-1}, Z_i || \gamma_i, \ell_i + s) & \text{when } i > 1. \end{cases} \end{aligned} \quad (5)$$

Note that  $\Pi$  and  $\Pi'$  can be same but we don't fix them here for generality.

**THEOREM D.1 (CCS IMPLIES prf-prng).** *Let  $\Pi$  and  $\Pi'$  be two CCS-secure schemes w.r.t. distributions  $\Sigma = (Z, C_Z)$  and  $\Sigma = (\Pi(Z, \gamma, \ell), C_Z || \gamma || \ell)$ , respectively then the PRF-PRNG  $\Pi_{pp}$  defined using  $\Pi$  and  $\Pi'$  as shown above will be prf-prng-secure under same input sources. More concretely, for all adversaries  $\mathcal{A}$  who make total  $q$   $\Pi_{pp}$  queries, we have*

$$\text{Adv}_{\Pi_{pp}[\Pi, \Pi']}^{\text{prf-prng}}(\mathcal{A}) \leq \text{Adv}_{\Pi}^{\text{CCS}}(\mathcal{B}) + \text{Adv}_{\Pi'}^{\text{CCS}}(\mathcal{C})$$

for some adversaries  $\mathcal{B}$  and  $\mathcal{C}$  making at most  $q_1$  and  $q_2$  queries to  $\Pi$  and  $\Pi'$ , respectively such that  $q_1 + q_2 = q$  and running in time given by the running time of  $\mathcal{A}$  plus  $\alpha_0 \cdot q$  for some constant  $\alpha_0$ .

**PROOF OF THEOREM D.1.** Let us recall from Sec. 7 and Fig. 8 that the CCS security of a KDF( $Z, \gamma, \ell$ ) w.r.t. the input source  $(Z, C_Z)$  implies that the output of a KDF query where the input has either unique  $\gamma$  value or independently sampled  $Z$  value is indistinguishable from  $\text{RF}_{Z, \gamma}[1 \dots \ell]$  i.e. uniform random binary string of length  $\ell$  bits. Let us now slightly abuse the notation for simplicity and denote by  $f$  and  $f'$  two functions that take inputs of the form  $(Z, \gamma, \ell)$  and return  $\text{RF}_{Z, \gamma} || \Pi[1 \dots \ell]$  and  $\text{RF}_{Z, \gamma} || \Pi'[1 \dots \ell]$  as outputs, respectively. This gives us

$$\text{Adv}_{\Pi_{pp}[\Pi, \Pi']}^{\text{prf-prng}}(\mathcal{A}) \leq \text{Adv}_{\Pi}^{\text{CCS}}(\mathcal{B}) + \text{Adv}_{\Pi'}^{\text{CCS}}(\mathcal{C}) + \text{Adv}_{\Pi_{pp}[f, f']}^{\text{prf-prng}}(\mathcal{A}).$$

Note that for any input  $(\sigma, I = Z || \gamma, \ell)$ ,  $\Pi_{pp}[f, f']$  as per Eqn. 5 will always return independently sampled uniform random strings  $(\sigma' || R)$  when  $I$  is unique (which implies that all returned chall-prf outputs in prf-prng games [4, Fig. 7] are indistinguishable from random strings). Similarly, for any input  $(\sigma, I = Z || \gamma, \ell)$ ,  $\Pi_{pp}[f, f']$  will always return independently sampled uniform random strings  $R$  when  $\sigma$  is the uncompromised current state (which implies that all returned chall-prng outputs in prf-prng games [4, Fig. 7] are indistinguishable from random strings). This implies  $\text{Adv}_{\Pi_{pp}[f, f']}^{\text{prf-prng}}(\mathcal{A}) = 0$  and hence the result.  $\square$

Clearly for Signal, setting  $\Pi = \text{Skye}$  (referring to KDF1 calls) and  $\Pi' = \text{FExp}$  (referring to KDF2 and KDF3 calls; which is CCS-secure as the input samples contain the current state value which is uniformly random and secret and thus can be used as the key to the PRF) gives us a  $\Pi_{pp}$  that covers all three types of KDF calls in Signal. The security here can be deduced from Theorem D.1 that says for all adversaries  $\mathcal{A}$  making a total of  $q$  queries to  $\Pi_{pp}$ , there exists some adversary  $\mathcal{B}$  making at most  $q$  queries and running in time given by the running time of  $\mathcal{A}$  plus some constant  $\alpha \cdot q$  such that

$$\text{Adv}_{\Pi_{pp}[\text{Skye}, \text{FExp}]}^{\text{prf-prng}}(\mathcal{A}) \leq \text{Adv}_{\text{Skye}}^{\text{CCS}}(\mathcal{B}).$$

We emphasize that as motivated in Sec. 9, this idea of using FExp in place of full Skye for KDF2 and KDF3 calls in Signal gives significant performance benefits.

## E PERFORMANCE DETAILS

In this section, we provide the benchmark tables Table 4 and 5 that correspond to the performance plots of Fig. 7a and 7b, respectively.

$n$	HKDF	Skye	HKDF	Skye	HKDF	Skye
	with AES-NI with SHA-NI		without AES-NI without SHA-NI		with AES-NI without SHA-NI	
1	2.22 ± 0.02	1.38 ± 0.01	6.51 ± 0.32	3.44 ± 0.03	6.47 ± 0.05	2.32 ± 0.03
2	4.48 ± 0.03	2.79 ± 0.02	12.97 ± 0.09	6.88 ± 0.04	12.95 ± 0.1	4.65 ± 0.03
3	6.7 ± 0.04	4.13 ± 0.03	19.43 ± 0.09	10.33 ± 0.22	19.4 ± 0.1	6.97 ± 0.04
4	8.92 ± 0.04	5.5 ± 0.08	25.93 ± 0.08	13.81 ± 0.49	25.88 ± 0.21	9.31 ± 0.06
5	11.14 ± 0.07	6.94 ± 0.05	32.48 ± 0.23	17.22 ± 0.21	32.36 ± 0.24	11.63 ± 0.09
6	13.34 ± 0.12	8.25 ± 0.05	38.82 ± 0.1	20.68 ± 0.05	38.86 ± 0.26	13.92 ± 0.27
7	15.68 ± 0.09	9.6 ± 0.08	45.42 ± 0.13	24.12 ± 0.46	45.33 ± 0.32	16.25 ± 0.12
8	17.98 ± 0.13	10.98 ± 0.1	51.82 ± 0.19	27.49 ± 0.07	51.86 ± 0.18	18.59 ± 0.14
9	20.07 ± 0.12	12.36 ± 0.12	58.38 ± 0.22	30.9 ± 0.27	58.35 ± 0.34	20.89 ± 0.07
10	22.33 ± 0.07	13.76 ± 0.31	64.89 ± 0.22	34.33 ± 0.29	64.81 ± 0.23	23.17 ± 0.12
20	44.44 ± 0.15	27.66 ± 1.38	129.82 ± 0.43	68.73 ± 0.18	129.48 ± 0.97	46.33 ± 0.18
30	66.68 ± 0.19	41.24 ± 0.28	194.63 ± 0.48	103.02 ± 0.58	194.16 ± 0.78	69.37 ± 0.43
40	88.66 ± 0.75	55.99 ± 0.25	260.27 ± 0.79	137.27 ± 0.68	258.92 ± 0.82	92.74 ± 0.31
50	111.6 ± 0.38	68.66 ± 0.53	324.22 ± 1.16	171.74 ± 1.59	323.76 ± 3.69	115.77 ± 0.71
60	133.79 ± 0.47	83.87 ± 0.51	389.55 ± 0.65	206.8 ± 0.52	388.39 ± 3.2	138.9 ± 1.01
70	155.49 ± 1.12	96.82 ± 0.58	453.72 ± 0.89	240.95 ± 0.64	453.18 ± 1.42	162.21 ± 1.07
80	178.02 ± 1.46	110.22 ± 0.81	518.02 ± 1.93	275.16 ± 1.84	517.8 ± 3.09	185.4 ± 0.66
90	201.36 ± 0.84	124.52 ± 0.47	582.89 ± 1.93	309.23 ± 1.77	582.69 ± 1.21	208.51 ± 1.31
100	222.8 ± 1.59	137.57 ± 1.05	649.14 ± 2.18	343.14 ± 1.99	646.73 ± 7.1	232.95 ± 10.46

**Table 4: The mean time ± standard deviation (in  $\mu$ s) required to encrypt  $n$  messages sent by one party to the other.**

$n$	HKDF	Skye	HKDF	Skye	HKDF	Skye
	with AES-NI with SHA-NI		without AES-NI without SHA-NI		with AES-NI without SHA-NI	
1	200.65 ± 1.12	194.51 ± 1.01	225.74 ± 1.86	204.82 ± 1.57	226.35 ± 1.37	201.42 ± 0.51
2	212.71 ± 1.81	203.19 ± 0.74	254.0 ± 0.8	221.24 ± 0.49	254.8 ± 2.01	214.08 ± 1.63
3	224.6 ± 1.79	211.19 ± 1.71	283.92 ± 7.51	238.06 ± 0.62	283.95 ± 1.05	226.12 ± 2.16
4	236.46 ± 0.69	219.88 ± 1.68	312.14 ± 0.9	254.73 ± 1.19	313.06 ± 1.02	238.73 ± 1.81
5	249.48 ± 1.52	229.85 ± 1.61	342.48 ± 2.38	272.02 ± 0.69	342.95 ± 1.1	251.37 ± 0.6
6	261.41 ± 1.79	237.91 ± 0.77	371.38 ± 1.48	288.83 ± 0.96	372.77 ± 8.9	264.32 ± 1.9
7	273.36 ± 1.4	246.89 ± 2.11	400.75 ± 1.27	305.93 ± 2.02	401.4 ± 1.54	276.52 ± 1.01
8	285.22 ± 1.13	255.93 ± 2.28	429.51 ± 1.06	322.96 ± 0.68	430.53 ± 1.11	288.91 ± 7.07
9	296.93 ± 1.6	264.1 ± 3.08	459.52 ± 1.13	337.91 ± 1.05	460.04 ± 3.28	300.99 ± 0.77
10	308.73 ± 0.64	272.63 ± 1.63	488.46 ± 1.15	355.45 ± 1.33	488.52 ± 1.71	313.48 ± 0.73
20	428.98 ± 0.89	361.04 ± 2.63	781.15 ± 6.48	522.62 ± 1.46	780.11 ± 2.08	437.46 ± 1.35
30	545.92 ± 1.08	446.9 ± 2.44	1071.67 ± 3.4	688.91 ± 3.74	1071.74 ± 9.22	559.56 ± 1.3
40	667.88 ± 8.06	535.84 ± 3.14	1364.68 ± 2.87	856.07 ± 2.04	1363.87 ± 8.33	684.44 ± 2.26
50	786.88 ± 1.78	621.98 ± 3.3	1652.43 ± 2.88	1021.6 ± 2.34	1655.0 ± 10.11	807.02 ± 2.37
60	905.83 ± 4.57	710.14 ± 10.74	1948.63 ± 12.21	1190.16 ± 4.19	1946.28 ± 3.94	935.55 ± 6.16
70	1025.61 ± 3.97	796.48 ± 5.8	2237.93 ± 4.99	1357.24 ± 2.92	2235.78 ± 4.38	1056.3 ± 2.32
80	1146.73 ± 2.48	888.52 ± 3.11	2528.96 ± 6.65	1523.47 ± 10.06	2530.07 ± 15.57	1184.09 ± 3.27
90	1263.38 ± 6.85	975.91 ± 6.12	2820.42 ± 9.19	1690.03 ± 6.13	2818.72 ± 9.55	1308.08 ± 3.49
100	1382.64 ± 2.92	1063.36 ± 8.75	3115.64 ± 18.99	1855.34 ± 3.72	3109.68 ± 17.06	1433.14 ± 7.3

**Table 5: The mean time ± standard deviation (in  $\mu$ s) required to encrypt (and decrypt)  $n$  messages that are sent (and received) by one party to (and from, respectively) the other.**