# Leftover Hash Lemma, Revisited

Boaz Barak \* Yevgeniy Dodis<sup>†</sup> Hugo Krawczyk<sup>‡</sup> Olivier Pereira<sup>§</sup>

Krzysztof Pietrzak<sup>¶</sup> Francois-Xa

Francois-Xavier Standaert

Yu Yu\*\*

September 3, 2011

#### Abstract

The famous *Leftover Hash Lemma* (LHL) states that (almost) universal hash functions are good randomness extractors. Despite its numerous applications, LHL-based extractors suffer from the following two limitations:

- Large Entropy Loss: to extract v bits from distribution X of min-entropy m which are  $\varepsilon$ -close to uniform, one must set  $v \leq m 2\log(1/\varepsilon)$ , meaning that the entropy loss  $L \stackrel{\text{def}}{=} m v \geq 2\log(1/\varepsilon)$ . For many applications, such entropy loss is too large.
- Large Seed Length: the seed length n of (almost) universal hash function required by the LHL must be at least  $n \ge \min(u v, v + 2\log(1/\varepsilon)) O(1)$ , where u is the length of the source, and must grow with the number of extracted bits.

Quite surprisingly, we show that both limitations of the LHL — large entropy loss and large seed — can be overcome (or, at least, mitigated) in various important scenarios. First, we show that entropy loss could be reduced to  $L = \log (1/\varepsilon)$  for the setting of deriving secret keys for a wide range of cryptographic applications. Specifically, the security of these schemes with an LHL-derived key gracefully degrades from  $\varepsilon$  to at most  $\varepsilon + \sqrt{\varepsilon 2^{-L}}$ . (Notice that, unlike standard LHL, this bound is meaningful even when one extracts more bits than the min-entropy we have!) Based on these results we build a general *computational extractor* that enjoys low entropy loss and can be used to instantiate a generic key derivation function for *any* cryptographic application.

Second, we study the soundness of the natural *expand-then-extract* approach, where one uses a *pseudorandom generator* (PRG) to expand a short "input seed" S into a longer "output seed" S', and then use the resulting S' as the seed required by the LHL (or, more generally, by any randomness extractor). We show that, in general, the expand-then-extract approach is not sound if the Decisional Diffie-Hellman assumption is true. Despite that, we show that it is sound either: (1) when extracting a "small" (logarithmic in the security of the PRG) number of bits; or (2) in minicrypt. Implication (2) suggests that the expand-then-extract approach is likely secure when used with "practical" PRGs, despite lacking a reductionist proof of security!

Finally, we combine our main results to give a very *simple and efficient* AES-based extractor, which easily supports *variable-length* messages, and is likely to offer our *improved entropy loss* bounds for any computationally-secure application, despite having a *fixed-length* seed.

<sup>\*</sup>Microsoft Research New England. Email: boaz@microsoft.com.

<sup>&</sup>lt;sup>†</sup>Department of Computer Science, New York University. Email: dodis@cs.nyu.edu.

<sup>&</sup>lt;sup>‡</sup>IBM Research. Email: hugo@ee.technion.ac.il.

<sup>&</sup>lt;sup>§</sup>Université Catholique de Louvain. Email: Olivier.Pereira@uclouvain.be.

<sup>¶</sup>CWI Amsterdam. Email: pietrzak@cwi.nl.

<sup>&</sup>lt;sup>||</sup>Université Catholique de Louvain. Email: fstandae@uclouvain.be.

<sup>\*\*</sup>Department of Computer Science, East China Normal University. Email: yuyu@yuyu.hk.

### 1 Introduction

The famous Leftover Hash Lemma [19] (LHL; see also [19] for earlier formulations) has found a huge number of applications in many areas of cryptography and complexity theory. In its simplest form, it states that universal hash functions [6] are good (strong) randomness extractors [33]. Specifically, if X is a distribution of min-entropy m over some space  $\mathcal{X}$ ,  $\mathcal{H}$  is a family of universal functions (see Definition 2.2) from  $\mathcal{X}$  to  $\{0,1\}^v$ , and H is a random member of  $\mathcal{H}$ , then, even conditioned on the "seed" H, the statistical distance between H(X) and the uniform distribution  $U_v$  on  $\{0,1\}^v$ is bounded by  $\sqrt{2^{-L}}$ , where  $L \stackrel{\text{def}}{=} m - v$ . The parameter L is defined as the entropy loss and it measures the amount of min-entropy "sacrificed" in order to achieve good randomness extraction. Thus, no application can tell apart the "extracted" randomness H(X) from uniform randomness  $U_v$ , with advantage greater than  $\varepsilon \stackrel{\text{def}}{=} \sqrt{2^{-L}}$ , even if the seed H is published (as long as H is independent of X).

The LHL is extremely attractive for many reasons. First, and foremost, it leads to simple and efficient randomness extractors, and can be used in many applications requiring good secret randomness. One such major setting is that of cryptographic key derivation, which is needed in many situations, such as privacy amplification [3], Diffie-Hellman key exchange [15, 27], biometrics [10, 4] and random number generators from physical sources of randomness [2, 1]. Second, many simple functions, such as the inner product or, more generally, matrix-vector multiplication, are universal. Such elegant functions have nice algebraic properties which can be used for other reasons beyond randomness extraction (for a few examples, see [28, 9, 31]). Third, many simple and efficient constructions of (almost) universal hash functions are known [6, 40, 32, 25], making LHL-based extractors the most efficient extractors to date. Finally, LHL achieves the optimal value of the entropy loss L = m - v sufficient to achieve the desired statistical distance  $\varepsilon$ . Specifically, LHL achieves  $L = 2 \log (1/\varepsilon)$ , which is known to be the the smallest possible entropy loss for any extractor [36].

Despite these extremely attractive properties, LHL-based extractors are not necessarily applicable or sufficient in various situations. This is primarily due to the following two limitations of the LHL: large entropy loss and large seed.

LARGE ENTROPY LOSS. In theory, the entropy loss of  $2 \log (1/\varepsilon)$  might appear quite insignificant, especially in the asymptotic sense. However, in practical situations it often becomes a deal-breaker, especially when applied to the setting of key derivation. In this case the main question is to determine the smallest min-entropy value m sufficient to extract a v-bit key with security  $\varepsilon$ . Minimizing this value m, which we call *startup entropy*, is often of critical importance, especially in entropy constrained scenarios, such as Diffie-Hellman key exchange (especially on elliptic curves)<sup>1</sup> or biometrics. For example, for the Diffie-Hellman key exchange, the value m corresponds to the size of the elliptic curve group, which directly affects efficiency. This is one of the reasons why statistical extractors are often replaced in practice with heuristic constructions based on cryptographic hash functions. See Section 3.4 for more on this point.

LARGE SEED. Another significant hurdle in the use of LHL comes from the fact that universal hash functions require long description, which means that LHL-based extractors have long seeds. Indeed, Stinson [40] showed that (perfectly) universal hash functions require the length of the seed to be linear in the length of the source X. More generally, even "good-enough" *almost* universal hash

<sup>&</sup>lt;sup>1</sup>Notice, in this application the exact source distribution is known, and, in principle, *deterministic* randomness extraction might be possible. Interestingly, state-of-the-art deterministic extractors [7] achieve noticeably *worse* security for a given entropy loss (see Theorems 8 and 14 of [7]) than the LHL-based (probabilistic) extractors. Thus, as of now, LHL remains the method of choice for this very important application scenario.

functions for LHL require seeds of length at least  $\min(|X| - v, v + 2\log(1/\varepsilon)) - O(1)$  [40], and, thus, must grow with the number of extracted bits. This large seed length makes it inconvenient in many applications of extractors (e.g., [5, 37, 27]), including any use of extractors for derandomization, where one must be able to enumerate over all the seeds efficiently.

Large (and variable-length) seeds are also inconvenient for standardized cryptographic applications where fixed-size keys, independent of the size of inputs, are favored (as in the case of block ciphers or cryptographic hash functions). When extractors are used in cryptographic settings, seeds are viewed as keys and hence fixed-size seeds are very desirable. In applications of extractors, where the attacker is assumed to be sufficiently limited as to not make the source X dependent on the seed (e.g., when extracting keys from biometrics, physical measurements or in the Diffie-Hellman key exchange), one might consider fixing a good *public* seed, and using it repeatedly with a fast provably secure extractor. As said, this is not possible with universal hash functions as their seed length must grow with the length of X.<sup>2</sup>

OUR RESULTS. Quite surprisingly, we show that both limitations of the LHL — large entropy loss and large seed — can be overcome or, at least, mitigated in various important scenarios. We describe these results below.

#### 1.1 Reducing the Entropy Loss

At first, reducing the entropy loss L might seem impossible since we already mentioned that any extractor must have entropy loss  $L \geq 2\log(1/\varepsilon) - O(1)$  [36]. However, the impossibility is for general applications of extractors, where we must ensure that the extracted string R cannot be distinguished from random by any statistical test D. In contrast, when extractors are used to derive cryptographic keys, we only care about limited types of statistical tests D. Concretely, the tests that correspond to the security game between the attacker A and the challenger C. For example, when deriving the key for a signature scheme, the only tests we care about correspond to the attacker seeing several signatures and then outputting a new signature. Namely, we only care that the probability of a successful forgery does not suddenly become non-negligible when the secret key is obtained using an extractor instead of being random. And since the signature scheme is assumed to be secure with a truly random key, we can restrict our attention to a very restricted class of statistical tests which almost never output 1. Similar restrictions on the distinguisher naturally arise for other cryptographic primitives, which gives us hope that the lower bound of [36] might be overcome in such settings.

GENERALIZED LHL AND APPLICATIONS. Indeed, we derive a tighter form of the LHL, called generalized LHL (see Theorem 3.1), which non-trivially depends on the type of distinguisher Dwe care about. Our improved bound contains a novel term informally measuring the standard deviation of the distinguisher's advantage (the standard LHL is a particular case where this term is bounded by 1). Applying this new bound to the analysis of cryptographic functions, we obtain much tighter bounds for the security of a wide class cryptographic applications. These include key derivation for all "unpredictability" applications, such as signatures, MACs, one-way functions, identification schemes, etc. More surprisingly, they also include key derivation for some prominent "indistinguishability" applications that include all stateless encryption schemes, both CPA- and CCA-secure and in the public- and symmetric-key settings, as well as weak pseudorandom functions. Specifically, in each of these cases, denote by  $\varepsilon$  the security of the cryptographic primitive (i.e.,

<sup>&</sup>lt;sup>2</sup>In theory one can build (non-LHL-based) extractors where the length of the seed H is roughly logarithmic in the length of the source X (see [17, 38] and many references therein). However, the resulting constructions are mainly of theoretical value and lose the extreme simplicity and efficiency of LHL-based extractors.

the best success probability or advantage of an attacker with certain resources) when keyed with a perfectly random v-bit key, and by  $\varepsilon'$  the corresponding security value when the key is derived from an imperfect m-bit entropy source via the LHL. We show (recall that L = m - v represents the entropy loss):

$$\varepsilon' \le \varepsilon + \sqrt{\varepsilon 2^{v-m}} = \varepsilon + \sqrt{\varepsilon 2^{-L}} \tag{1}$$

COMPARING WITH STANDARD LHL. Let us first compare this bound with the regular  $\varepsilon + \sqrt{2^{-L}}$ LHL bound. The latter required  $L \geq 2\log(1/\varepsilon)$  to achieve the same type of security  $O(\varepsilon)$  as with the ideal randomness. Using our improved bound, we show that only half of that amount,  $L = \log(1/\varepsilon)$ , already suffices. In fact, not only do we get improved bounds on the entropy loss L, but we also get meaningful security bounds for arbitrary values of L, even negative ones (when the entropy loss becomes "entropy gain")! E.g., standard LHL does not give anything for  $L \leq 0$ , while we achieve significant  $\varepsilon' \approx \sqrt{\varepsilon}$  security for L = 0 (no entropy loss!), and even start to "gracefully borrow" security from our application when we extract more bits than the min-entropy of the source, up to  $L = -\log(1/\varepsilon)$  (i.e.,  $v = m + \log(1/\varepsilon)$ ).

COMPUTATIONAL EXTRACTOR WITH IMPROVED LOSS. Although our improved bound, as stated, is not applicable to all cryptographic applications (the most important omission being pseudorandom functions and stream ciphers), in Section 3.3 we use our results to build *general-purpose* key derivation function for *any* (computationally-secure) cryptographic application, while providing the full entropy benefits derived from Equation (1). The scheme combines any LHL-based extractor with any (weak) pseudorandom function family.

### 1.2 Reducing the Seed Length

EXPAND-THEN-EXTRACT. A natural idea to reduce the seed length is to use a *pseudorandom* generator (PRG) to expand a short "input seed" S into a longer "output seed" S', and then use the resulting S' as the seed required by the LHL, or, more generally, by any randomness extractor. Let us call this natural approach *expand-then-extract*. Of course, as long as one hides the short Sand uses the long S' as the public seed for the extractor, the extracted bits are pseudorandom. But is it possible to ensure the pseudorandomness of the extractor's output if the actual short seed S is made *public*? Had this been the case, we would obtain efficient LHL-based extractors with a short seed and, moreover, an extractor whose seed length is independent of the length of the input, as desired for the practical scenarios discussed earlier.

COUNTER-EXAMPLE. In Section 4.1 we show that the expand-then-extract approach will not work in general. We construct a simple PRG (which is secure under the Decisional Diffie-Hellman (DDH) assumption) and an extractor (which is a natural, perfectly universal hash function), where the output of the extractor — on any (e.g., even uniform) distribution — can be efficiently distinguished from random with probability close to 1, when given the short seed S used to generate the pseudorandom long seed S' for the extractor. Despite the above, we also show two positive results which nicely complement our counter-example.

EXTRACTING FEW BITS. First, in Section 4.2 we show that the expand-then-extract approach always works provided the number of extracted bits v is "small". Here "small" means logarithmic in the security level of the PRG, which could range from  $O(\log k)$  to  $\Omega(k)$  (where k is the security parameter), depending on whether PRG is assumed to be polynomially or exponentially hard. Quite interestingly, in this case we do not even have to settle for pseudorandom bits: our small number v of extracted bits is actually statistically random, as long as the PRG is secure against circuits whose size is exponential in v. The intuition for this result comes from the fact that we can test, in time exponential in v, whether a given *n*-bit extractor seed s' is "good" or "bad" for our source X. We also know that most *random* long seeds  $s' \leftarrow U_n$  must be good. Hence, by the PRG security, the same must be true for "most" *pseudorandom* seeds  $s' \leftarrow Prg(U_k)$ , which is precisely what we need to show.

SECURITY IN minicrypt. Second, although our original counterexample is fairly simple and natural, it involves an assumption (DDH) from the "public-key world". In Section 4.3, we show, somewhat surprisingly, that such "public-key" type assumption is indeed *necessary* for any counter-example. We do this by showing that the expand-then-extract approach is sound in minicrypt [23] (i.e. in a hypothetical world where pseudorandom generators exist, but public-key cryptography does not). In particular, we construct a simple 2-message protocol (built from a PRG and an extractor) which constitutes a secure key-agreement (KA) protocol for any PRG/extractor combination for which the expand-then-extract approach is *insecure*. Since our protocol only has 2 messages, we even get semantically secure public-key encryption (PKE). In fact, the protocol has a very special form which even implies a 2-round oblivious transfer (OT) protocol secure against an honest-but-curious receiver. Hence, since no such KA/PKE/OT exist in minicrypt, expand-then-extract must be secure.

The protocol (whose slight generalization to handle side information is depicted in Figure 1) proceeds as follows. Alice samples a random seed S and sends  $S' = \Pr(S)$  to Bob. Bob samples a source X and a random bit b. If b = 0, Bob sends Alice the output  $\operatorname{Ext}(X; S)$  of the extractor; otherwise he sends a random string R. Alice, who knows S, can recover b (by assumption). In contrast, Eve, who only sees  $S' = \Pr(S)$ , cannot tell apart  $\operatorname{Ext}(X; S)$  from R (and thus guess b) as this would mean she can distinguish  $\operatorname{Prg}(S)$  from uniform.

PRACTICAL INTERPRETATION. This leads to the following practical interpretation of our results indicating that using the expand-then-extract approach with common pseudorandom primitives, such as AES, is secure in spite of a lack of direct (reductionist) proof of security. Indeed, consider the expand-then-extract scheme implemented via AES (in some stream cipher mode). Our results show that, if this extraction scheme fails, then we have found a public-key encryption (and even oblivious transfer) scheme that is as secure as the AES as a block cipher! Moreover, the resulting PKE has a very restrictive form, where the secret key is a PRG seed S, and the public-key is the PRG output  $S' = \Pr g(S)$ . (E.g., in the case of AES, the public key is simply the evaluation of AES on several distinct points.) As we argue in Section 4.3, the existence of such a PKE appears to be extremely unlikely, and would be a major breakthrough given current state-of-the-art. Thus, our results give strong evidence that the expand-then-extract approach might be secure "in practice", — when used with "fast" ciphers (like AES), — despite being (generally) insecure "in theory"!

EXTENSIONS. We also remark that all our results elegantly handle side information Z the attacker might have about our source X (as advocated by [10], such "average-case" extractors are very handy in cryptographic applications), and also generalize to the case of *almost* universal hash functions. Finally, one can combine our results regarding improved entropy loss in Section 3 with those regarding the soundness of the expand-then-extract approach in Section 4, to obtain an efficient (computationally secure) LHL-based extractor with short seed and improved startup entropy. In Section 5, we give a simple example of such an extractor, by using the AES block cipher and the inner-product universal hash function. See Equation (25).

#### 1.3 Related Work

Hast [18] also observed that for certain cryptographic applications, the relevant attackers correspond to restricted classes of distinguishers, which allowed him to obtain improved security bounds when the Goldreich-Levin hardcore bit [16] is used as a "computational" randomness extractor. This result is incomparable to ours. On the one hand, we consider general (multi-bit) LHL-based extractors and not just the single bit inner-product function (which is the form of the Goldreich-Levin predicate). On the other hand, Hast was working in the computational setting, and had to make an explicit reduction from the distinguisher to the predictor of the source X, which is not required in our setting.

We also mentioned the notion of *slice extractors* defined by Radhakrishnan and Ta-Shma [36], which limits the type of statistical tests to "rare distinguishers". To the best of our understanding, this definition was not motivated by applications, but rather was a convenient "parametrization" on a road to other results. Still, this notion roughly correspond to the setting of key derivation for authentication applications, when the attacker rarely succeeds. Interestingly, [36] showed a lower bound for the entropy loss of slice extractors (which was lower than that of general extractors), and matched this lower bound by an existential construction. As it turns out, our improved LHL immediately gives a *constructive* way to match this lower bound, showing that LHL-based extractors are optimal slice extractors in terms of the entropy loss. For completeness, this connection is sketched in Appendix A.

In a very different (non-cryptographic) context of building hash tables, Mitzenmacher and Vahdan [29] also observed that improved bounds on the "entropy loss" could be obtained when the standard deviation of the "distinguisher" is much less than 1. In their setting the entropy loss was the minimum entropy required from the input stream to hash well, and the distinguisher was the characteristic function of a set of occupied buckets.

We note that our "win-win" result for the expand-then-extract approach is similar in spirit to several other "win-win" results [13, 34, 11, 35], where a (hypothetical) adversary for one task is turned into a "surprisingly useful" protocol for a seemingly unrelated task. Among the above, the result most similar to ours is [13], where a PRG is used to expand the key for "forward secure storage", which is a concept related to "locally computable" extractors.

On the more practical side of our results, particularly in what refers to key derivation, it is worth mentioning the work of [8, 27] that analyze constructions of key derivation functions (KDFs) based on cryptographic hash functions. These constructions do not use standard, generic assumptions, such as pseudorandomness, but build on specific modes of operations on their compression function f, and rely on dedicated, sometimes idealized, assumptions. Under such idealized assumptions, these schemes support situations where the KDF needs to be modeled as a random oracle, or where the source only has "unpredictability entropy" [22]. On the other hand, outside of the random oracle heuristics, much of the analysis of [8, 27] studied sufficient conditions on compression function f and/or the source input distribution, under which cryptographic hash functions are "universal enough" so as to apply the standard LHL. As such, these analyses suffer the same drawbacks as any other LHL-based extractor. In particular, our results regarding the improved entropy loss for LHL-based extractor should carry over to improve the results of [8, 27], while our results on the expand-then-extract approach could be viewed as partial justification of the heuristic where a fixed-description-length compression function is replaced by random in most (but not all) of the analyses of [8, 27].

# 2 Standard Leftover Hash Lemma

NOTATION. For a set S, we let  $U_S$  denote the uniform distribution over S. For an integer  $v \in \mathbb{N}$ , we let  $U_v$  denote the uniform distribution over  $\{0, 1\}^v$ , the bit-strings of length v. For a distribution or random variable X we write  $x \leftarrow X$  to denote the operation of sampling a random x according to X. For a set S, we write  $s \leftarrow S$  as shorthand for  $s \leftarrow U_S$ .

MIN-ENTROPY AND EXTRACTORS. The min-entropy of a random variable X is defined as  $\mathbf{H}_{\infty}(X) \stackrel{\text{def}}{=} -\log(\max_{x} \Pr[X = x])$ . In cryptographic applications, one often uses the average min-entropy of a random variable X conditioned on another random variable Z. This is defined as

$$\widetilde{\mathbf{H}}_{\infty}(X|Z) \stackrel{\text{def}}{=} -\log\left(\mathbb{E}_{z\leftarrow Z} \left[\max_{x} \Pr[X=x|Z=z]\right]\right) = -\log\left(\mathbb{E}_{z\leftarrow Z} \left[2^{-\mathbf{H}_{\infty}(X|Z=z)}\right]\right)$$

where  $\mathbb{E}_{z \leftarrow Z}$  denotes the expected value over  $z \leftarrow Z$ , and measures the worst-case predictability of X by an adversary that may observe a correlated variable Z.

We denote with  $\Delta_D(X, Y)$  the advantage of a circuit D in distinguishing the random variables  $X, Y: \Delta_D(X, Y) \stackrel{\text{def}}{=} | \Pr[D(X) = 1] - \Pr[D(Y) = 1] |$ . The statistical distance between two random variables X, Y is defined by

$$\mathsf{SD}(X,Y) \stackrel{\text{def}}{=} \frac{1}{2} \sum_{x} |\Pr[X=x] - \Pr[Y=x]| = \max_{D} \quad \Delta_D(X,Y)$$

where the maximum is taken over all (potentially computationally unbounded) D. Given side information Z, we write  $\Delta_D(X, Y|Z)$  and SD(X, Y|Z) as shorthands for  $\Delta_D((X, Z), (Y, Z))$  and SD((X, Z), (Y, Z)), respectively.<sup>3</sup>

An extractor [33] can be used to extract uniform randomness out of a weakly-random value which is only assumed to have sufficient min-entropy. Our definition follows that of [10], which is defined in terms of conditional min-entropy.

**Definition 2.1 (Extractors)** An efficient function  $\text{Ext} : \mathcal{X} \times \{0,1\}^n \to \{0,1\}^v$  is an (averagecase, strong)  $(m,\varepsilon)$ -extractor (for space  $\mathcal{X}$ ), if for all X, Z such that X is distributed over  $\mathcal{X}$  and  $\widetilde{\mathbf{H}}_{\infty}(X|Z) \geq m$ , we get

$$\mathsf{SD}(\mathsf{Ext}(X;S), U_v \mid (S,Z)) \leq \varepsilon$$

where  $S \equiv U_n$  denotes the coins of Ext (called the seed). The value L = m - v is called the entropy loss of Ext, and the value n is called the seed length of Ext.

UNIVERSAL HASHING AND LEFTOVER HASH LEMMA. We now recall the definition of universalhashing [6, 40] and the leftover-hash lemma [19], which states that universal hash functions are also good extractors.

**Definition 2.2** ( $\rho$ -Universal Hashing) A family  $\mathcal{H}$  of (deterministic) functions h :  $\mathcal{X} \to \{0,1\}^v$  is a called  $\rho$ -universal hash family (on space  $\mathcal{X}$ ), if for any  $x_1 \neq x_2 \in \mathcal{X}$  we have  $\Pr_{h \leftarrow \mathcal{H}}[h(x_1) = h(x_2)] \leq \rho$ . When  $\rho = 1/2^v$ , we say that  $\mathcal{H}$  is universal.

We can finally state the Leftover Hash Lemma (LHL). (Multiple versions of this lemma have appeared; we use the formulation of [41, Theorem 8.1], augmented by [10, Lemma 2.4] for the conditional entropy case; see [19] and references therein for earlier formulations.)

**Lemma 2.1 (Leftover-Hash Lemma)** Assume that the family  $\mathcal{H}$  of functions  $h : \mathcal{X} \to \{0,1\}^v$ is a  $\frac{1+\gamma}{2^v}$ -universal hash family. Then the extractor  $\mathsf{Ext}(x;h) \stackrel{\text{def}}{=} h(x)$ , where h is uniform over  $\mathcal{H}$ , is an  $(m,\varepsilon)$ -extractor, where  $\varepsilon = \frac{1}{2} \cdot \sqrt{\gamma + \frac{2^v}{2^m}} = \frac{1}{2} \cdot \sqrt{\gamma + \frac{1}{2^L}}$  (recall, L = m - v is the entropy loss). In particular,  $\frac{1+3\varepsilon^2}{2^v}$ -universal hash functions yield  $(v + 2\log(1/\varepsilon), \varepsilon)$ -extractors with entropy loss  $L = 2\log(1/\varepsilon)$ .

<sup>&</sup>lt;sup>3</sup>Notice,  $\Delta_D(X, Y|Z) \leq \mathbb{E}_{z \leftarrow Z}[\Delta_D(X|_{Z=z}, Y|_{Z=z})]$ , but  $\mathsf{SD}(X, Y|Z) = \max_D \mathbb{E}_{z \leftarrow Z}[\Delta_D(X|_{Z=z}, Y|_{Z=z})]$ .

## **3** Reducing the Entropy Loss

As we mentioned, the entropy loss of  $2 \log (1/\varepsilon)$  is optimal when one is concerned with general distinguishers D [36]. As we show, in various cryptographic scenarios we only care about a somewhat restrictive class of distinguishers, which will allow us to reduce the entropy loss for such applications.

In Section 3.1, we start with a generalization of the LHL (Theorem 3.1), which will include a novel term measuring the *standard deviation* of the distinguisher's advantage, and then derive some useful special cases (Theorem 3.2). In Section 3.2, we then apply our tighter bound to derive improved entropy loss bounds for various cryptographic applications, including bounds for so called "strongly secure" applications (Theorem 3.3), all authentication applications (Theorem 3.4) and some privacy applications (Theorem 3.5), including chosen plaintext secure encryption (Theorem 3.6) and weak PRFs (Theorem 3.7). In Section 3.3 we further extend our results to get a generic key derivation function with improved entropy loss for *any* computationally secure application, including stream ciphers and PRFs. In Section 3.4 we compare our bounds to the heuristic bounds obtained using the random oracle model.

#### 3.1 Generalized LHL

COLLISION PROBABILITY AND *c*-VARIANCE. Given a distribution *Y*, its collision probability is  $\operatorname{Col}(Y) \stackrel{\text{def}}{=} \sum_{y} \Pr[Y = y]^2 \leq 2^{-\mathbf{H}_{\infty}(Y)}$ . Given a joint distribution (Y, Z), we let  $\operatorname{Col}(Y|Z) \stackrel{\text{def}}{=} \mathbb{E}_{z}[\operatorname{Col}(Y|Z = z)] \leq 2^{-\widetilde{\mathbf{H}}_{\infty}(Y|Z)}$ . We also use the notation  $(U_{\mathcal{Y}}, Z)$  to denote the probability distribution which first samples  $(y, z) \leftarrow (Y, Z)$ , but then replaces *y* by an independent, uniform sample from  $U_{\mathcal{Y}}$ . Finally, given a random variable *W* and a constant *c*, we define its *c*-variance as  $\operatorname{Var}_{c}[W] \stackrel{\text{def}}{=} \mathbb{E}[(W - c)^{2}]$  and its *c*-standard deviation as  $\sigma_{c}[W] \stackrel{\text{def}}{=} \sqrt{\operatorname{Var}_{c}[W]}$ . When  $c = \mathbb{E}[W]$ , we recover the standard notion of variance and standard deviation,  $\operatorname{Var}[W]$  and  $\sigma[W]$ , and notice that these are the smallest possible values for  $\operatorname{Var}_{c}[W]$  and  $\sigma_{c}[W]$ . Still, we will later find it easier to use (slightly weaker) bounds obtained with specific values of  $c \in \{0, \frac{1}{2}\}$ .

We start with an useful lemma of independent interest which generalizes Lemma 4.5.1 of [29].

**Lemma 3.1** Assume (Y, Z) is a pair of correlated random variables distribution on a set  $\mathcal{Y} \times \mathcal{Z}$ . Then for any (deterministic) real-valued function  $f: \mathcal{Y} \times \mathcal{Z} \to \mathbb{R}$  and any constant c, we have

$$|\mathbb{E}[f(Y,Z)] - \mathbb{E}[f(U_{\mathcal{Y}},Z)]| \leq \sigma_c[f(U_{\mathcal{Y}},Z)] \cdot \sqrt{|\mathcal{Y}|\mathsf{Col}(Y|Z) - 1}$$
(2)

*Proof.* If suffices to prove the claim for c = 0, by using g(y, z) = f(y, z) - c in place of f(y, z), as such a change does not affect the left hand side.

Denote  $q_z = \Pr[Z = z], p_{y|z} = \Pr[Y = y \mid Z = z]$  and also recall the Cauchy-Schwartz inequality  $|\sum a_i b_i| \leq \sqrt{(\sum a_i^2) \cdot (\sum b_i^2)}$ . We have

$$\begin{aligned} |\mathbb{E}[f(Y,Z)] - \mathbb{E}[f(U_{\mathcal{Y}},Z)]| &= \left| \sum_{y,z} q_z \cdot f(y,z) \cdot p_{y|z} - \sum_{y,z} q_z \cdot f(y,z) \cdot \frac{1}{|\mathcal{Y}|} \right| \\ &= \left| \sum_{y,z} \left( \sqrt{\frac{q_z}{|\mathcal{Y}|}} \cdot f(y,z) \right) \cdot \left( \sqrt{|\mathcal{Y}|q_z} \cdot \left( p_{y|z} - \frac{1}{|\mathcal{Y}|} \right) \right) \right| \\ &\leq \sqrt{\left( \sum_{y,z} \frac{q_z}{|\mathcal{Y}|} \cdot f(y,z)^2 \right) \cdot \left( |\mathcal{Y}| \sum_{y,z} q_z \cdot \left( p_{y|z} - \frac{1}{|\mathcal{Y}|} \right)^2 \right)} \\ &= \sqrt{\mathbb{E}[f(U_{\mathcal{Y}},Z)^2] \cdot |\mathcal{Y}| \cdot \left( \sum_{y,z} q_z \cdot p_{y|z}^2 - \frac{1}{|\mathcal{Y}|} \right)} \\ &= \sigma_0[f(U_{\mathcal{Y}},Z)] \cdot \sqrt{|\mathcal{Y}|\text{Col}(Y|Z) - 1} \end{aligned}$$

The useful feature of the bound given in Equation (2) comes from the fact that the value  $\sigma_c[f(U_{\mathcal{Y}}, Z)]$ does not depend on the actual distribution Y used to replace the uniform distribution, while the value  $\sqrt{|\mathcal{Y}| \text{Col}(Y|Z) - 1}$  does not depend on the function f whose average we are trying to preserve (by using Y in place of  $U_{\mathcal{Y}}$ ). In particular, we get the following corollary which will allow us to eventually get all our improved bounds.

**Theorem 3.1 (Generalized LHL)** Let (X, Z) be some joint distribution over  $\mathcal{X} \times \mathcal{Z}$ ,  $\mathcal{H} = \{h : \mathcal{X} \to \{0,1\}^v\}$  be a family of  $\frac{1+\gamma}{2^v}$ -universal hash functions, H be a random member of  $\mathcal{H}$ , and let  $L \stackrel{\text{def}}{=} \widetilde{\mathbf{H}}_{\infty}(X|Z) - v$  be the entropy loss. Then, for any constant  $c \in [0,1]$  and any (possibly probabilistic) distinguisher D(r,h,z), we have

$$\Delta_D(H(X), U_v \mid (H, Z)) \le \sigma_c \left[ \Pr_D[D(U_v, H, Z) = 1] \right] \cdot \sqrt{\gamma + \frac{1}{2^L}}$$
(3)

*Proof.* We use Lemma 3.1 with  $\mathcal{Y} = \{0,1\}^v \times \mathcal{H}, Y = (H(X),H)$  and  $f(y,z) = f(r,h,z) \stackrel{\text{def}}{=} \Pr_D[D(r,h,z) = 1]$ , where the probability is taken only over the coins of D (so that  $f(U_v, H, Z)$  is a random variable, with range [0,1], over the choice of  $U_v, H, Z$ ). Comparing the needed bound in Equation (3) with the given bound in Equation (2), and recalling that  $L = \widetilde{\mathbf{H}}_{\infty}(X|Z) - v$ , it suffices to show that

$$\mathsf{Col}((H(X),H) \mid Z) \le \frac{1+\gamma+2^{v} \cdot 2^{-\mathbf{H}_{\infty}(X\mid Z)}}{2^{v} \cdot |\mathcal{H}|} = \frac{1}{|\mathcal{H}|} \cdot \left(\frac{1+\gamma}{2^{v}} + 2^{-\tilde{\mathbf{H}}_{\infty}(X\mid Z)}\right) \tag{4}$$

Fix any value Z = z and let  $X_z$  be the conditional distribution on X when Z = z. Then, the value  $Col(H(X_z), H)$  can be interpreted as a probability that  $(H(X_z), H) = (H'(X'_z), H')$ , where  $H, H', X_z, X'_z$  are four independent samples of the appropriate distributions. We get:

$$\begin{aligned} \mathsf{Col}((H(X_{z}), H)) &= & \Pr_{H, H', X_{z}, X'_{z}}[(H(X_{z}), H) = (H'(X'_{z}), H')] \\ &= & \Pr_{H, H'}[H = H'] \cdot \Pr_{H, X_{z}, X'_{z}}[H(X_{z}) = H(X'_{z})] \\ &\leq & \frac{1}{|\mathcal{H}|} \cdot \left( \Pr_{H, X_{z}, X'_{z}}[H(X_{z}) = H(X'_{z}) \mid X_{z} \neq X'_{z}] + \Pr_{X_{z}, X'_{z}}[X_{z} = X'_{z}] \right) \\ &\leq & \frac{1}{|\mathcal{H}|} \cdot \left( \frac{1 + \gamma}{2^{v}} + 2^{-\mathbf{H}_{\infty}(X_{z})} \right) \end{aligned}$$

where the last inequality follows from the universality of  $\mathcal{H}$  and the definition of min-entropy. Taking expectation over  $z \leftarrow Z$ , and recalling that  $2^{-\widetilde{\mathbf{H}}_{\infty}(X|Z)} = \mathbb{E}_{z \leftarrow Z}[2^{-\mathbf{H}_{\infty}(X|Z=z)}]$ , we get the desired Equation (4).

Equation (3) bounds the advantage of D in distinguishing real and extracted randomness using two terms. The second term (under the square root) depends on the universality of  $\mathcal{H}$  and the entropy loss L (but not on D). The novel term is the c-standard deviation  $\sigma_c[\Pr_D[D(U_v, H, Z) = 1]]$  of D, which we will simply call c-standard deviation of D and denote  $\sigma(D, c)$ . Intuitively, it measures how "concentrated" the expected output of D is to some value c in the "ideal setting" (when fed  $U_v$  rather than H(X)). We notice that for any D and any  $c \in [0, 1]$ , the c-standard deviation  $\sigma(D, c) \leq 1$ . Plugging this trivial bound in Equation (3) removes the dependence on D, and (essentially)<sup>4</sup> gives us the statement of the standard LHL from Lemma 2.1. As mentioned, though, this forces the entropy loss to be at least  $2\log(1/\varepsilon)$  to achieve security  $\varepsilon$ . Below, we show several special cases when we can upper bound  $\sigma(D, c)$  by roughly  $\sqrt{\varepsilon}$ , which means that the entropy loss L only needs to be roughly  $\log(1/\varepsilon)$  (say, with perfectly universal  $\mathcal{H}$ ) to achieve security  $\varepsilon$ .

**Theorem 3.2** Let (X, Z) be some joint distribution over  $\mathcal{X} \times \mathcal{Z}$ ,  $\mathcal{H} = \{h : \mathcal{X} \to \{0, 1\}^v\}$  be a family of  $\frac{1+\gamma}{2^v}$ -universal hash functions, H be a random member of  $\mathcal{H}$ ,  $L \stackrel{\text{def}}{=} \widetilde{\mathbf{H}}_{\infty}(X|Z) - v$  be the entropy loss, and D(r, h, z) be some (possibly probabilistic) distinguisher. Then, for each of the values  $\varepsilon$  defined in scenarios (a)-(c) below it holds:

$$\Delta_D(H(X), U_v \mid (H, Z)) \le \sqrt{\varepsilon \cdot \left(\gamma + \frac{1}{2^L}\right)}$$
(5)

(a) Assume for some  $c, \delta, \tau \in [0, 1]$ ,  $\varepsilon = \tau^2 + \delta$  and the following condition is satisfied:<sup>5</sup>

$$\Pr_{r \leftarrow U_v, h \leftarrow \mathcal{H}, z \leftarrow Z} \left[ \left| \Pr_D[D(r, h, z) = 1] - c \right| \ge \tau \right] \le \delta$$
(6)

- (b) Assume  $\Pr[D(U_v, H, Z) = 1] \le \varepsilon$  (where probability is taken over  $U_v, H, Z$  and the coins of D).
- (c) For fixed r, h, and z, define the distinguisher D'(r, h, z) as follows. First, make two independent samples  $\tilde{d}, d \leftarrow D(r, h, z)$ . Then, if  $\tilde{d} = 1$ , return d else return (1 d). Assume further that  $\Pr[D'(U_v, H, Z) = 1] \leq \frac{1}{2} + 2\epsilon$ .

Proof. Each part will follow from Equation (3) by showing that  $\operatorname{Var}_c[\operatorname{Pr}_D[D(U_v, H, Z) = 1]] \leq \varepsilon$  for some  $c \in [0, 1]$ . To simplify the notation, let q = (r, h, z),  $f(q) = \operatorname{Pr}_D[D(q) = 1]$  and  $Q = (U_v, H, Z)$ . Recalling also the definition of c-variance, we need to show that for some  $c \in [0, 1]$ ,

$$\mathbb{E}_{q \leftarrow Q}[(f(q) - c)^2] \le \varepsilon \tag{7}$$

Proof of part (a). With our notation, condition (a) states that  $\Pr_{q \leftarrow Q}[(f(q) - c)^2] \ge \tau^2] \le \delta$ . Since also  $(f(q) - c)^2 \le 1$  for all q, we get  $\mathbb{E}_{q \leftarrow Q}[(f(q) - c)^2] \le \tau^2 \cdot (1 - \delta) + 1 \cdot \delta \le \varepsilon$ .

Proof of part (b). Here we set c = 0 and note that  $\mathbb{E}_{q \leftarrow Q}[(f(q) - c)^2] = \mathbb{E}_{q \leftarrow Q}[f(q)^2] \leq \mathbb{E}_{q \leftarrow Q}[f(q)] \leq \varepsilon$ , where  $f(q)^2 \leq f(q)$  because  $0 \leq f(q) \leq 1$ .

<sup>&</sup>lt;sup>4</sup>The exact bound claimed in Lemma 2.1 follows when  $\sigma(D,c) \leq \frac{1}{2}$ , which is true for c = 1/2.

<sup>&</sup>lt;sup>5</sup>Note that the condition below implies  $|\Pr[D(U_v, H, Z) = 1] - c| \le \tau + \delta$ .

Proof of part (c). Here we set  $c = \frac{1}{2}$ , so it suffices to show that  $\mathbb{E}_{q \leftarrow Q}[(f(q) - \frac{1}{2})^2] \leq \varepsilon$ . Recalling the definition of f and that by assumption  $2\varepsilon \geq \mathbb{E}_{q \leftarrow Q}[\Pr[D'(q) = 1] - \frac{1}{2}]$ , it suffices to show that for all q we have

$$\Pr[D'(q) = 1] - \frac{1}{2} \ge 2 \cdot \left(\Pr[D(q) = 1] - \frac{1}{2}\right)^2$$

Letting  $a = \Pr[D(q) = 1]$  and  $b = \Pr[D'(q) = 1]$ , this is equivalent to the claim that  $b \ge 2(a - \frac{1}{2})^2 + \frac{1}{2} = a^2 + (1 - a)^2$ . The latter, however, is immediate from the definition of D'(q), which outputs 1 only if both samples of D(q) are the same.

In the following section, we demonstrate the use of Theorem 3.2 by concentrating on the important case of *key derivation* using LHL, where the value  $\varepsilon$  will essentially correspond to the "cryptographic security" of the application at hand.

#### 3.2 Improved LHL for Key Derivation

Consider any cryptographic primitive P (e.g., signature, encryption, etc.), which uses randomness  $R \in \{0,1\}^v$  to derive its secret (and, public, if needed) key(s). Without loss of generality, we can assume that R itself is the secret key. In the "ideal" setting,  $R = U_v$  is perfectly uniform and independent from whatever side information Z available to the attacker. In the "real setting", the key owner has a randomness source, represented by a random variable X and possibly correlated with the attacker's side information Z. It then samples a universal hash function H using (fresh) public randomness and uses the extracted value R = H(X) as its key. We would like to argue that if P is " $\varepsilon$ -secure" in the ideal setting (against attackers with resources<sup>6</sup> less than T), then P is also " $\varepsilon'$ -secure" in the real setting (against attackers with resources less than  $T' \approx T$ ), where  $\varepsilon'$  is not much larger than  $\varepsilon$ . Of course, to have a hope of achieving this,  $\mathcal{H}$  must be "universal-enough" and  $L = \widetilde{\mathbf{H}}_{\infty}(X|Z) - v$  must "high-enough". To parameterize this, we will sometimes explicitly write  $(L, \gamma)$ -real model to denote the real model above, where  $\mathcal{H}$  is  $(1 + \gamma)2^{-v}$ -universal and  $\widetilde{\mathbf{H}}_{\infty}(X|Z) \geq v + L$ . We formalize this general setting as follows.

ABSTRACT SECURITY GAMES. We assume that the security of P is defined via an interactive game between a probabilistic attacker A(h, z) and a probabilistic challenger C(r). Here one should think of h and z as particular values of the hash function and the side information, respectively, and r as a particular value used by the challenger in the key generation algorithm of P. We note that C only uses the secret key r and does not depend on h and z. In the ideal setting, where  $r \leftarrow U_v$ , the attacker A does not use the values h and z (and anyway the optimal values of h and z can be hardwired into A in the non-uniform model), yet, for notation convenience, we will still pass h and z to A even in the ideal setting.

At the end of the game, C(r) outputs a bit b, where b = 1 indicates that the attacker "won the game". Since C is fixed by the definition of P (e.g., C runs the unforgeability game for signature or the semantic security game for encryption, etc.), we denote by  $D_A(r, h, z)$  the (abstract) distinguisher which simulates the entire game between A(h, z) and C(r) and outputs the bit b, and by  $\operatorname{Win}_A(r, h, z) = \Pr[D_A(r, h, z) = 1]$  the probability that A(h, z) wins the game against C(r). With this notation, the probability of winning the game in the "real setting" is given by the random variable  $\operatorname{Win}_A(H(X), H, Z)$ , and the same probability in the ideal setting becomes  $\operatorname{Win}_A(U_v, H, Z)$ . Moreover, the difference between these probabilities is simply the distinguishing advantage of  $D_A$ 

<sup>&</sup>lt;sup>6</sup>We use the word "resource" to include all the efficiency measures we might care about, such as running time, circuit size, number of oracle queries, etc.

of telling apart real and extracted randomness when given H, Z:

$$|Win_A(H(X), H, Z) - Win_A(U_v, H, Z)| = \Delta_{D_A}(H(X), U_v \mid (H, Z))$$
(8)

As we justify next, to argue the security of P in the real setting assuming its security in the ideal setting, it is sufficient for us to argue that the above distinguishing advantage is "small" for all legal attackers A. And since the security of P will usually restrict the power of attackers A (hence, also the power of abstract distinguishers  $D_A$ ), we may use the results of Theorem 3.2 to argue better bounds on the entropy loss  $L = \widetilde{\mathbf{H}}_{\infty}(X|Z) - v$ .

**Definition 3.1** Let c = 0 for unpredictability applications P (signature, MAC, one-way function, etc.) and  $c = \frac{1}{2}$  for indistinguishability applications P (encryption, pseudorandom function/permutation, etc.). Fix also the  $(1 + \gamma)2^{-v}$ -universal hash family  $\mathcal{H}$  and the joint distribution (X, Z) satisfying  $\widetilde{\mathbf{H}}_{\infty}(X|Z) \geq v + L$ , so that the real and the ideal model are well-defined.

We say that P is  $(T, \varepsilon)$ -secure in the ideal model if for all attackers A with resources less than T, we have  $\operatorname{Win}_A(U_v, H, Z) \leq c + \varepsilon$ .

Similarly, P is  $(T', \varepsilon')$ -secure in the real model if for all attackers A have resources less than T', we have  $\operatorname{Win}_A(H(X), H, Z) \leq c + \varepsilon'$ .

Triangle inequality coupled with Equation (8) immediately yields the following Corollary.

**Lemma 3.2** Fix L and  $\gamma$  defining the real and the ideal models. Assume P is  $(T, \varepsilon)$ -secure in the ideal model, and for all attackers A with resources less than T' (where  $T' \leq T$ ) we have  $\Delta_{D_A}(H(X), U_v \mid (H, Z)) \leq \delta$ . Then P is  $(T', \varepsilon + \delta)$ -secure in the  $(L, \gamma)$ -real model.

We are now ready to apply Lemma 3.2 and Theorem 3.2 to various cryptographic primitives P. Below, we let  $c \in \{0, \frac{1}{2}\}$  be the constant governing the security of P (0 for unpredictability and 1/2 for indistinguishability applications).

#### 3.2.1 General Bound for "Strongly Secure" Applications

Fix h and z and consider any attacker A. We say that a key  $r \in \{0,1\}^v$  is  $\tau$ -bad w.r.t. A(h,z) if  $|Win_A(r,h,z) - c| \geq \tau$ . We then say that P is strongly  $(T,\tau,\delta)$ -secure in the ideal model if for all h and z and all attackers A(h,z) with resources less than T, we have

$$\Pr_{\tau \leftarrow U_v}[r \text{ is } \tau \text{-bad w.r.t. } A(h, z)] \le \delta$$
(9)

Namely, the fraction of keys r on which A has advantage more than  $\tau$  is at most  $\delta$ . It is easy to see that strong  $(T, \tau, \delta)$ -security implies (regular)  $(T, \tau + \delta)$ -security, since even on  $\tau$ -bad keys the advantage of A is at most 1. On the other hand, strong  $(T, \tau, \delta)$ -security w.r.t. A clearly implies that the abstract distinguisher  $D_A$  satisfies Equation (6). Therefore, we can apply part (a) of Theorem 3.2 to Lemma 3.2 and get:

**Theorem 3.3** Fix L and  $\gamma$  defining the real and the ideal models. Assume P is strongly  $(T, \tau, \delta)$ -secure in the ideal model. Then P is  $(T, \varepsilon)$ -secure in the ideal model and  $(T, \varepsilon')$ -secure in the  $(L, \gamma)$ -real model, where

$$\varepsilon \le \tau + \delta$$
 and  $\varepsilon' \le \tau + \delta + \sqrt{(\tau^2 + \delta) \cdot \left(\gamma + \frac{1}{2^L}\right)}$  (10)

In particular, when  $\tau = \delta = \varepsilon/2$ , we get  $\varepsilon' \leq \varepsilon + \sqrt{\varepsilon(\gamma + 2^{-L})}$ .

DISCUSSION. Strong security assumes that any attacker really fails on all but a negligible fraction of keys. Although many cryptographic primitives achieve this level of security, a notable exception is the one-time pad. Of course, this is not surprising, since we really do not expect to get any better security bounds for the one-time pad than what is possible with the standard LHL. However, it is instructive to see what goes wrong. Consider an attacker A who claims to distinguish a message  $m_0$  whose first bit is 0 from a message  $m_1$  whose first bit is 1. When getting the ciphertext  $e = m_b \oplus r$ , A outputs the first bit of e (call it b'). Clearly, the overall advantage of A is zero, as  $\Pr[b = b'] = \Pr[r_1 = 0] = 1/2$ , where  $r_1$  is the first bit of our key r. However, A is correct with probability 1 on all keys r whose first bit is 0, and incorrect on all keys r whose first bit is 1. Thus, as expected, one-time pad is not strongly secure according to our definition, since on every key r the adversary's odds of success are bounded away from 1/2. A similar problem happens with other "deterministic indistinguishability" primitives, such as pseudorandom generators, functions and permutations.

On the other hand, Theorem 3.3 readily applies to unforgeability applications (corresponding to c = 0). For example, one can easily show that regular  $(T, \tau \delta)$ -security implies strong  $(T, \tau, \delta)$ -security. This simple observation can already give improved bounds on the entropy loss for this setting. However, instead of pursuing this (rather loose) implication, in Section 3.2.2 we will directly prove a much stronger bound: see Theorem 3.4.

Similarly, it appears reasonable to assume that concrete "probabilistic indistinguishability" primitives, such as many semantically secure encryption schemes, satisfy a decent enough kind of strong security. Once again, though, we will directly prove good security bounds for such applications in Section 3.2.3, without making assumptions about their strong security.

To summarize, the bounds arising from strong security are useful and interesting, but better bounds can often be obtained by direct analysis, as we show in Section 3.2.2 and Section 3.2.3.

#### 3.2.2 Improved Bound for Unpredictability Applications

Recall, authentication applications correspond to c = 0, and include signature schemes, MACs, oneway functions/permutations, etc. In this case  $(T, \varepsilon)$ -security in the ideal model implies that for any T-bounded attacker A,  $\mathbb{E}[\operatorname{Win}_A(U_v, H, Z)] \leq \varepsilon$ . Recalling the definition of the abstract distinguisher  $D_A$ , this is the same as  $\Pr[D_A(U_v, H, Z) = 1] \leq \varepsilon$ , which is precisely the pre-condition for part (b) of Theorem 3.2. Thus, combining Equation (5) with Lemma 3.2, we immediate get:

**Theorem 3.4** Fix L and  $\gamma$  defining the real and the ideal models. Assume authentication primitive P (corresponding to c = 0) is  $(T, \varepsilon)$ -secure in the ideal model. Then P is  $(T, \varepsilon')$ -secure in the  $(L, \gamma)$ -real model, where

$$\varepsilon' \le \varepsilon + \sqrt{\varepsilon \left(\gamma + \frac{1}{2^L}\right)}$$
 (11)

In particular, if  $\gamma = 0$  and  $L = \log(1/\varepsilon)$ , then  $\varepsilon' \leq 2\varepsilon$ . Moreover, when  $\gamma = 0$ , the security bound is meaningful even for negative entropy "loss"  $0 \geq L > -\log(1/\varepsilon)$ , when one extracts more bits than the min-entropy  $\widetilde{\mathbf{H}}_{\infty}(X|Z)$  and "borrows the security deficit" from the ideal security of P.

Intuitively, Theorem 3.4 uses the fact that for authentication applications one only cares about distinguishers which almost never output 1, since the attacker almost never forges successfully.

#### 3.2.3 Improved Bound for Some Indistinguishability Applications

We now move to the more difficult case of indistinguishability applications, where c = 1/2. In general, we do not expect to outperform the standard LHL, as illustrated by the one-time pad

example. Quite surprisingly, we show that the for a class of applications, including stateless chosen plaintext attack (CPA) secure encryption, one can still get improved bounds as compared to the standard LHL. Specifically, as long as the primitive P allows the attacker A to "test" its success before playing the actual "challenge" from C, we still get significant savings. We start by defining the general type of security games where our technique applies.

BIT-GUESSING GAMES. As usual the game is played by the attacker A = A(h, z) and the challenger C(r). The game can have an arbitrary structure, except the winning condition is determined as follows. At some point A asks C for a "challenge". C flips a random bit  $b \in \{0, 1\}$ , and sends A a value e = e(b, r). The game continues in an arbitrary way and ends with A making a guess b'. A wins if b = b'.

So far, this still includes all standard indistinguishability games, including the one-time pad. The extra assumption we make is the following. For any valid attacker A having resources less than T' there exists another valid attacker A' (having somewhat larger resources  $T \ge T'$ ) such that:

- (1) The execution between A' and C(r) defines four bits  $b, b', \tilde{b}, \tilde{b}'$ , such that the joint distribution of  $(b, b', \tilde{b}, \tilde{b}')$  is the same as two *independent* tuples (b, b') obtained when A runs with C(r).
- (2) The bits b and b' are precisely the secret bit of C and the guess of A', so that A' wins iff b = b'.
- (3) A' learns if  $\tilde{b} = \tilde{b}'$  before outputting b'.

We will call such indistinguishability games (T', T)-simulatable.

**Theorem 3.5** Fix L and  $\gamma$  defining the real and the ideal models. Assume indistinguishability primitive P (corresponding to c = 1/2) is (T', T)-simulatable and  $(T, \varepsilon)$ -secure in the ideal model. Then P is  $(T', \varepsilon')$ -secure in the  $(L, \gamma)$ -real model, where

$$\varepsilon' \le \varepsilon + \sqrt{\varepsilon \left(\gamma + \frac{1}{2^L}\right)} \tag{12}$$

In particular, if  $\gamma = 0$  and  $L = \log(1/\varepsilon)$ ,  $\varepsilon' \leq 2\varepsilon$ . Moreover, when  $\gamma = 0$ , the security bound is meaningful even for negative entropy "loss"  $0 \geq L > -\log(1/\varepsilon)$ , when one extracts more bits than the min-entropy  $\widetilde{\mathbf{H}}_{\infty}(X|Z)$  and "borrows the security deficit" from the ideal security of P.

*Proof.* [(of Theorem 3.5)] By Lemma 3.2, it is sufficient to show that for all attackers A with resources less than T' we have

$$\Delta_{D_A}(H(X), U_v \mid (H, Z)) \le \sqrt{\varepsilon(\gamma + 2^{-L})}$$
(13)

In turn, using part (c) of Theorem 3.2, (13) will follow if we can show that

$$\Pr[D'_A[U_v, H, Z] = 1] \le \frac{1}{2} + 2\varepsilon \tag{14}$$

where  $D'_A(r, h, z)$  runs  $D_A$  twice (with independent coins), and then reverses the second answer if the first answer is 0.

Let us examine the outputs of  $D_A$  and  $D'_A$  more closely.  $D_A$  outputs a bit d which is 1 iff b = b'(namely, A won the game).  $D'_A$  (independently) runs  $D_A$  twice, and observes the first answer  $\tilde{d}$ , which is 1 iff  $\tilde{b} = \tilde{b}'$ . If so,  $D'_A$  simply returns the second answer d of  $D_A$ , which is 1 if b = b'. Otherwise,  $D'_A$  reverses d, which, for the bit-guessing games, is equivalent to reversing the second output b' of A to 1 - b'.

As we can see, the behavior of  $D'_A$  is almost identical to the behavior of the legal *T*-bounded attacker A' assumed by the (T', T)-simulatability of P. Indeed, the run of A' defines two independent runs of A (by condition (1)). A' also learns whether the first run was successful (e.g., if  $\tilde{d} = 1$ ), by condition (3). The only difference is that A' always outputs the value b', which results in the final bit being d (by condition (2)). On the other hand,  $D'_A$  reverses its output to 1 - b' if the first run was not successful. So we can simply define another T-bounded<sup>7</sup> attacker  $\tilde{A}$  who simply flips the final output bit of A' if the first run was not successful, which is within the rules by condition (2). With this in mind, we get a T-bounded attacker  $\tilde{A}$  which wins the game precisely when  $D'_A$ outputs 1; in symbols,  $D_{\tilde{A}} \equiv D'_A$ .

Since  $\tilde{A}$  is a valid *T*-bounded attacker and  $D_{\tilde{A}} \equiv D'_A$ , the  $(T, \varepsilon)$ -security of *P* implies that  $\Pr[D'_A[U_v, H, Z] = 1] \leq \frac{1}{2} + \varepsilon$ , proving Equation (14).

Below we show that stateless CPA-secure (public- or symmetric-key) encryption schemes are simulatable, where, as expected, the "resources" T are roughly doubled compared to T'. We start with the slightly simpler public-key setting. In this case the attacker gets the public-key pk, and can completely simulate the run of the challenger C, including the winning condition, since C does not use the secret key r. In particular, for any CPA-attacker A we define A' which first runs A against "virtual C", using the public-key to see if A correctly guessed the encrypted message. After that A'can simply honestly simulate the run of A against the "actual C". It is obvious that this simulation satisfies conditions (1)-(3), at the expense of doubling the running time of A' as compared with A. We can also easily adjust it to the setting of chosen *ciphertext* attack (CCA) security, where the challenger C can decrypt arbitrary ciphertexts for A. In this case, A still needs "actual C" to simulate the decryption queries during the first "virtual run". The only slight subtlety is to ensure that A does not ask to decrypt the second (actual) challenge ciphertext during the first (virtual) run, but the chances of that are easily seen to be negligible.

Things are similar for the symmetric-key setting, except A no longer gets the public-key, and needs the help of "actual C" to simulate both the encryption queries and the challenge query of A during the "virtual" run. Namely, to argue real CPA-security of P against q encryption queries, we must assume ideal CPA security of P against (2q + 1) encryption queries, where the first q + 1queries will be used to simulate the q encryption and 1 challenge queries of A during the first run. Notice, since  $2q + 1 \ge 1$ , we indeed do not cover the one-time pad (or more generally, one-time secure symmetric-key encryption), since its security does not allow any encryption queries.

Summarizing the above discussion, and using Theorem 3.5, we get:

**Theorem 3.6** Fix L and  $\gamma$  defining the real and the ideal models, and set  $\varepsilon' = \varepsilon + \sqrt{\varepsilon(\gamma + 2^{-L})}$ . Assume P is public-key encryption scheme which is  $\varepsilon$ -secure, in the ideal model, against attackers with running time  $2t + t_{enc}$ , where  $t_{enc}$  is the runtime of the encryption process. Then P is  $\varepsilon'$ -secure, in the  $(L, \gamma)$ -real model, against attackers with running time t.

Similarly, assume P is a stateless symmetric-key encryption scheme which is  $\varepsilon$ -secure, in the ideal model, against attackers with running time 2t + O(1) and making 2q + 1 encryption queries. Then P is  $\varepsilon'$ -secure, in the  $(L, \gamma)$ -real model, against attackers with running time t and making q encryption queries.

LIMITATIONS AND EXTENSIONS. Unfortunately, several other indistinguishability primitives, such as pseudorandom generators, functions or permutations, do not appear to be simulatable. The

<sup>&</sup>lt;sup>7</sup>We assume the decision to flip the output bit does not really change the resources of A'.

problem seems to be in verifying the winning condition (condition (3) of simulatability), since this has to be done with respect to the actual secret key r not known to the attacker. For PRFs (or PRPs), it is tempting to solve the problem by using an equivalent definition, where the attacker can learn the value of PRF at any point, but then, as a challenge, must distinguish the value of the PRF at an *un-queried* point from random. Although this variant allows the attacker to check the winning condition during the first "virtual" run, it now creates a different problem in that the challenge point during the second "actual" run might have been queried during the first run, making such an attacker A' invalid.

Interestingly, the above "fix" works for a useful relaxation of PRFs, known as *weak PRFs* (wPRFs). Here the attacker only gets to see the values of the PRF at several *random points*, and has to distinguish a new value from random. Assuming a large enough input domain, the probability of collision between the PRF values revealed in the first first run and challenged in the second run, is negligible, which allows to complete the (valid) simulation. Similarly, it works for a slightly stronger relaxation of PRFs, known as *random-challenge* PRFs. As with wPRFs, A gets as the challenge a real-or-random evaluations of the PRF at a random point, but can additionally query the PRF at arbitrary points different from the challenge point. In Section 3.3 we show that wPRFs are all we need to apply our results to a generic key derivation function.

#### **3.2.4** Numeric Examples

STRONG SECURITY. Assume X is a random element of 160-bit elliptic group, Z is empty,  $\mathcal{H}$  is (perfectly) universal, and we wish to extract 128-bit AES key R to be used in some block-cipher based MAC (e.g., CBC-MAC). Thus, L = 160 - 128 = 32. Using standard LHL, we can achieve security at most  $2^{-L/2} = 2^{-16}$ , irrespective of the MAC security. However, assume that (for some T) the MAC is strongly secure with  $\tau = 2^{-50}$  and  $\delta = 2^{-70}$  (i.e., the probability of forgery is at most  $2^{-50}$  for all but  $2^{-70}$  fraction of keys), which seems quite reasonable. Using Equation (10), we get security  $\varepsilon' \leq 2^{-50} + 2^{-70} + \sqrt{(2^{-100} + 2^{-70})2^{-32}} \approx \frac{3}{2} \cdot 2^{-50}$ . This bound is not only much better than  $2^{-16}$ , but is essentially the same as the regular  $2^{-50}$ -security we started from in the ideal setting!

UNPREDICTABILITY. Consider again the setting AES-based MAC, where the key size v = 128. Now, let us assume that (for some T and query bound  $q \leq 2^{24}$ ) the MAC is  $2^{-80}$ -secure, so  $\varepsilon = 2^{-80}$ . As before X is uniform over some elliptic curve group, but now let us consider the group size u as a parameter, and see how the our security bounds degrade with u. In all the cases, assume perfectly universal hash family. Using the standard LHL, we get  $\varepsilon' = \varepsilon + 2^{(v-u)/2} = 2^{-80} + 2^{64} \cdot 2^{-u/2}$ . This bound is meaningless when group size  $u \leq 128$ , gives very low  $\varepsilon' \approx 2^{-16}$  for u = 160, and needs u = 288 to have  $\varepsilon' = 2\varepsilon = 2^{-79}$ . Using Theorem 3.4, we get  $\varepsilon' = \varepsilon + \sqrt{\varepsilon 2^{v-u}} = 2^{-80} + 2^{24} \cdot 2^{-u/2}$ . For a 128-bit group, it gives a "respectable"  $\varepsilon' \approx 2^{-40}$ , while not losing any entropy. For u = 160, it gives already good  $\varepsilon' \approx 2^{-56}$ , and only needs the group size u = 208 to achieve  $\varepsilon' = 2\varepsilon = 2^{-79}$ . Remarkably, with our bounds one may get non-trivial results even when the group's entropy u is *smaller* than the required output, namely, *gaining entropy rather than losing it!*. For example, consider deriving an HMAC-SHA-256 key from a 160-bit group. Assuming HMAC-SHA-256 has security of  $2^{-180}$  (against an attacker of certain resources and with a 256-bit random key), our results show that deriving the key from a much shorter 160-bit group element still provides a non-trivial security bound of  $2^{-42}$  while gaining 96 bits of entropy.

INDISTINGUISHABILITY. Consider the AES-CBC encryption setting with keys of size v = 128. Let us assume that (for some t) the resulting CBC encryption is  $2^{-50}$ -secure (so  $\varepsilon = 2^{-50}$ ) against CPA attacker running in time  $2t + t_{enc}$  and making a million encryption queries. As before, X is uniform over some elliptic curve group and we will consider the group size u as a parameter. In all the cases, assume a perfectly universal hash family. Using the standard LHL, we get  $\varepsilon' = \varepsilon + 2^{(v-u)/2} = 2^{-50} + 2^{64} \cdot 2^{-u/2}$  (although do not need to lose a factor of 2 in the run-time and the number of encryption queries). This bound is meaningless when group size  $u \leq 128$ , gives very low  $\varepsilon' \approx 2^{-16}$  for u = 160, and needs u = 228 to have  $\varepsilon' = 2\varepsilon = 2^{-49}$ . Using Theorem 3.6, we get  $\varepsilon' = \varepsilon + \sqrt{\varepsilon 2^{v-u}} = 2^{-50} + 2^{39} \cdot 2^{-u/2}$ . This bound is non-trivial even for u = 100 bits, despite extracting 128 bits! For 128-bit group, it gives a decent  $\varepsilon' \approx 2^{-25}$ , while not losing any entropy. For u = 160, it gives already good  $\varepsilon' \approx 2^{-41}$ , and only needs the group size u = 188 to achieve  $\varepsilon' = 2\varepsilon = 2^{-49}$ .

Note that all of the above examples apply to deriving keys from Diffie-Hellman values  $g^{xy}$  over the corresponding groups assuming the Decisional Diffie-Hellman assumption and setting the conditioning variable Z to the pair  $(g^x, g^y)$ .

#### 3.3 A Generic Key Derivation Function

So far we have discussed the applications of our generalized Leftover Hash Lemma and Theorem 3.2 to the derivation of cryptographic keys in entropy-constraint environments for specific applications. Although our analysis covers many such applications, it does not cover all, including PRFs and stream ciphers. In this section we make a simple observation which allows us to overcome this limitation and design a *generic* key derivation function (KDFs) which is (computationally) secure for any cryptographic application while still enjoying the same entropy loss savings. The idea is to compose universal hash functions a *weak* PRF (wPRF), where the random input to the wPRF now becomes part of the extractor seed, and use the fact that wPRFs fall under the class of *simulatable* applications as defined in Section 3.2.3.

Specifically, we define the KDF Ext on the basis of a  $(1 + \gamma)/2^v$ -universal hash family  $\mathcal{H}$  with v-bit outputs and a wPRF F taking a k-bit input w and a v-bit key r. Without loss of generality, we assume the output length of  $F_r(w)$  is also v bits, as one can always turn an  $\ell$ -bit output wPRF F', where  $\ell < v$ , into a v-bit output wPRF F, by setting  $F_r(w_1, \ldots, w_t) \stackrel{\text{def}}{=} (F'_r(w_1) \ldots F'_r(w_t))$ , where  $t = \lceil v/\ell \rceil$ . The public seed s for the KDF Ext is a pair s = (h, w), where h is a random universal function from  $\mathcal{H}$  and w is a random element in the domain of F. We then define  $\text{Ext}(x, (h, w)) \stackrel{\text{def}}{=} F_{h(x)}(w)$ ; i.e., the initially extracted value h(x) is used as a wPRF key, which is then used to evaluate F on w.

**Theorem 3.7** Assume  $F, \mathcal{H}$ , Ext are defined as above, which fixes the values of v, k and  $\gamma$ . Let X be any distribution of min-entropy  $m \ge v + L$ , and assume that F is an  $\varepsilon$ -secure wPRF against attackers A with of size 2t learning one value of F before the (random) challenge:

$$|\Pr[A(W_1, W_2, F_R(W_1), F_R(W_2)) = 1] - \Pr[A(W_1, W_2, F_R(W_1), U_v) = 1]| \le \varepsilon$$

Then no attacker B of size t can distinguish the output of Ext from uniform with advantage better than  $\varepsilon' \stackrel{\text{def}}{=} \varepsilon + \sqrt{\varepsilon(\gamma + 2^{-L})} + \frac{1}{2^k}$ . Namely,

$$|\Pr[B(\mathsf{Ext}(X, (H, W)), H, W, Z)) = 1] - \Pr[B(U_v, H, W, Z) = 1]| \le \varepsilon'$$

In particular, if the extracted bits are used as a key for any application P having  $\delta$ -security against attackers of size  $t' \leq t$  in the ideal model, then P is  $\delta + \varepsilon'$  secure against the same class of attackers in the  $(L, \gamma)$ -real model.

*Proof.* Follows in the same manner as the proof of Theorem 3.6 for symmetric-key encryption with q = 0, with wPRF playing the role of symmetric-key encryption, since the pseudorandomness

of the extracted bits  $\text{Ext}(X, (H, W)) = F_{H(X)}(w)$  corresponds to wPRF security in the real model under no message attack (q = 0). The extra term  $1/2^k$  in  $\varepsilon'$  corresponds to the probability of collision between the simulated run and the real run.

We also notice that if one needs to extract multiple keys for several applications, we can simply use the output of our computational KDF as a seed of a regular PRG or PRF, since such applications are now "covered".

**Remark 3.1** For the case of deriving multiple keys, as above, we notice that the wPRF step is not needed provided all the keys are for cryptographic applications covered by our technique (i.e., strongly secure, unpredictable, or simulatable primitives). Namely, in such a case we can directly use the initially extracted key H(X) as a seed for the (regular) PRF/PRG to derive all the required keys. This allows us to avoid increasing the seed length by k bits, and saves one application of wPRF. The proof of this claim follows by a simple hybrid argument (which we omit). In general, though, the wPRF-based solution is preferable, as it adds considerable generality at a relatively moderate cost. Additionally, the wPRF security  $\varepsilon$  (against a single query) in Theorem 3.7 might be considerably lower than the security  $\delta$  of the application P for which we extract of key (e.g., if P is a stateless CPA-secure encryption scheme based on wPRFs). In this case, the error term  $\sqrt{\varepsilon(\gamma + 2^{-L})}$  in Theorem 3.7 might be much smaller than the corresponding term  $\sqrt{\delta(\gamma + 2^{-L})}$ , obtained by direct analysis using tools in Section 3.2.

#### **3.4** Comparison to Random Oracle

It is also instructive to compare our new generalized LHL (GLHL) bound with the heuristic "dream" bound obtained by using cryptographic hash functions H as extractors. In the latter case, we claim that, when extracting a v-bit key for a primitive P with security  $\varepsilon$  using a random oracle (RO), the security degrades by roughly  $\varepsilon 2^{-L}$ , where the "entropy loss" L is still defined to be m - v.<sup>8</sup> Indeed, in the RO model the value H(X) can be distinguished from random only if the attacker queried the random oracle of X. Since X has min-entropy m, this happens with probability at most  $T/2^m$ , where T is an upper bound on the adversary's running time. Also, since our our application P is  $\varepsilon$ -secure against attackers of complexity T, we must have  $\varepsilon \leq T/2^v$ , since the attacker should not be able to find the correct key by exhaustive search with probability more that  $\varepsilon$  (moreover, this bound might be tight for many practical applications, where v is equal to the security parameter). Thus,  $T \approx \varepsilon \cdot 2^v$ , and the attacker's chance of distinguishing H(X) from uniform is at most

$$\varepsilon_{\mathsf{RO}} = \frac{T}{2^m} \approx \frac{\varepsilon \cdot 2^v}{2^{v+L}} = \varepsilon 2^{-L}$$

In contrast, our GLHL analysis achieved security degradation  $\varepsilon_{\mathsf{GLHL}} = \sqrt{\varepsilon 2^{-L}}$  (recall also that traditional LHL analysis had degradation  $\varepsilon_{\mathsf{LHL}} = \sqrt{2^{-L}}$ ). Thus, not surprisingly, the heuristic random oracle bound is still better than the GLHL bound, as it does not involve the square root operation. Of course, it is achieved by resorting to an *idealized heuristic*, while the GLHL bound applies to very simple extractors in a *rigorous and provable* way. Moreover, our results show, for the first time, that the dependency on the security of the cryptographic primitive (represented by  $\varepsilon$ ) is not just the property of idealized functions but applies to any LHL-based extractor. In particular, just like the RO bound (and unlike standard LHL), our GLHL bound is meaningful even for the

<sup>&</sup>lt;sup>8</sup>Of course, one can always make the output length v arbitrarily large in the random oracle model, much like in the standard model one can (computationally) amplify v using a pseudorandom generator. So, for a meaningful comparison between LHL and RO, one should think of v as equal to security parameter in both settings.

entropy gain situations  $0 \ge L > -\log(1/\varepsilon)$ , where we "borrow" the security from the application at hand.

In practical terms, the RO heuristic allows min-entropy m = v to achieve comparable security  $\varepsilon_{\text{RO}} = O(\varepsilon)$ , GLHL analysis permits  $m = v + \log(1/\varepsilon)$  for  $\varepsilon_{\text{GLHL}} = O(\varepsilon)$ , while the standard LHL analysis requires  $m = v + 2\log(1/\varepsilon)$  for  $\varepsilon_{\text{LHL}} = O(\varepsilon)$ . Thus, while our analysis does not quite match the RO heuristic, it significantly narrows the gaps between provable and heuristic startup entropy bounds.

### 4 Reducing the Seed Length

In this section we study the soundness of the natural expand-then-extract approach described in the Introduction. After providing the needed definitions below, in Section 4.1 we show that, unfortunately, the expand-then-extract approach will *not* work. In general (under the DDH assumption). We show two positive results which nicely complement our counter-example. First, in Section 4.2 we show that the expand-then-extract approach always works (and even gives a regular, not just computational, extractor!) provided *the number of extracted bits v is "small*". Second, in Section 4.3 we show that the expand-then-extract approach is sound in minicrypt, which shows that (1) any counter-example must use assumptions from cryptomania (such as DDH); (2) the expand-then-extract approach is likely secure when used with "practical" PRGs, despite the lack of a reductionist proof.

NEGLIGIBLE AND FRIENDS. We use k to denote a security parameter. A function  $\mu : \mathbb{N} \to [0, 1]$  is *negligible* if for any c > 0 there is a  $k_0$  such that  $\mu(k) \leq 1/k^c$  for all  $k \geq k_0$ . To the contrary,  $\mu$  is *non-negligible* if for some c > 0 we have  $\mu(k) \geq 1/k^c$  for infinitely many k. Throughout,  $\operatorname{negl}(k)$  denotes a negligible function in k.

A function  $\tau(\cdot) : \mathbb{N} \to [0,1]$  is overwhelming if  $1 - \tau(\cdot)$  is negligible. A function  $\phi : \mathbb{N} \to [0,1]$  is noticeable if for some c > 0 there is an  $k_0$  such that  $\phi(k) \ge 1/k^c$  for all  $k \ge k_0$ . Note that non-negligible is not the same as noticeable. For example,  $\mu(k) \stackrel{\text{def}}{=} k \mod 2$  is non-negligible but not noticeable.

COMPUTATIONAL EXTRACTORS AND PRGS. Recall that with  $\Delta_D(X,Y) \stackrel{\text{def}}{=} | \Pr[D(X) = 1] - \Pr[D(Y) = 1] |$  we denote the advantage of a circuit D in distinguishing the random variables X and Y. Let  $\mathcal{D}_t$  denote the class of all probabilistic circuits of size t. With  $\mathsf{CD}_t(X,Y) = \max_D \Delta_D(X,Y)$  we denote the computational distance of X and Y, here the maximum is over  $D \in \mathcal{D}_t$ . When  $t = \infty$  in unbounded, we recover the notion of statistical distance  $\mathsf{SD}(X,Y)$ . When  $X = X_k$  and  $Y = Y_k$  are families of distributions indexed by the security parameter k, we will say that X and Y are computationally indistinguishable, denoted  $X \approx_c Y$ , if for every polynomial t(.),  $\mathsf{CD}_{t(k)}(X_k, Y_k) = \mathsf{negl}(k)$ .

**Definition 4.1 (Computational Extractor)** We say that an efficient function  $\mathsf{Ext} : \mathcal{X} \times \{0, 1\}^n \to \{0, 1\}^v$  is an (average-case, strong) computational *m*-extractor (for space  $\mathcal{X}$ ), if for all efficiently samplable X, Z such that X is distributed over  $\mathcal{X}$  and  $\widetilde{\mathbf{H}}_{\infty}(X|Z) \geq m$  (here  $\mathcal{X}, n, v, m$  are all indexed by a security parameter k)

$$(\mathsf{Ext}(X;S),S,Z) \approx_c (U_v,S,Z)$$

**Definition 4.2 (Pseudorandom Generator)** A length increasing function  $Prg : \{0,1\}^k \to \{0,1\}^n$ is a pseudorandom generator (PRG) if  $Prg(U_k) \approx_c U_n$ . We also say that Prg is  $(T, \delta)$ -secure if

<sup>&</sup>lt;sup>9</sup>In order to avoid an extra parameter, we simply assume wlog that the seed length of our PRG is equal to the security parameter k.

 $\mathsf{CD}_T(\mathsf{Prg}(U_k), U_n) \leq \delta.$ 

COMPUTATIONAL EXTRACTORS WITH SHORT SEEDS? We are ready to formalize our main question: Is it safe to expand an extractor seed using a PRG?

**Hypothesis 1** [Expand-then-Extract] If Ext is an  $(m(k), \varepsilon(k))$ -extractor with seed length n(k) where  $\varepsilon(.)$  is negligible and  $\operatorname{Prg} : \{0,1\}^k \to \{0,1\}^{n(k)}$  is a pseudorandom generator, then  $\operatorname{Ext'}$  defined as

$$\mathsf{Ext}'(x;s) = \mathsf{Ext}(x;\mathsf{Prg}(s))$$

is a computational m-extractor.

#### 4.1 Counter-Example: Expanding Seeds is Insecure in General

In this section we show that, unfortunately, Hypothesis 1 is wrong in general.

**Theorem 4.1 (Hypothesis 1 wrong assuming DDH)** Under the DDH assumption, there exists a pseudorandom generator Prg(.) and a strong extractor Ext(.;.) (which is a perfectly universal hash function) such that  $Ext'(x;s) \stackrel{\text{def}}{=} Ext(x; Prg(s))$  can be efficiently distinguished from uniform on any input distribution (i.e. Ext' is not a computational extractor.)

*Proof.* [of Theorem 4.1] Let  $\mathcal{G}$  be a prime order p cyclic group with generator g where the DDH problem is hard. Then  $\operatorname{Prg} : \mathbb{Z}_p^3 \to \mathcal{G}^6$  defined as

$$Prg(a, b, c) = (g^{a}, g^{b}, g^{ab}, g^{ac}, g^{bc}, g^{abc})$$
(15)

is a a secure pseudorandom generator [30]. Let  $\mathsf{Ext}:\mathbb{Z}_p^3\times\mathcal{G}^6\to\mathcal{G}^2$  be

$$\mathsf{Ext}((x, y, z); (A, B, C, D, E, F)) = (A^x B^y C^z, D^x E^y F^z)$$

It is easy to see that Ext is a perfectly universal hash function from  $\mathbb{Z}_p^3 \to \mathcal{G}^2$  (and, by Lemma 2.1, strong  $(2\log p + 2\log(1/\varepsilon), \varepsilon)$ -extractor). Now consider the distribution

$$[\mathsf{Ext}((x, y, z); \mathsf{Prg}(a, b, c)) , (a, b, c)] = [(g^{ax}g^{by}g^{abz}, g^{acx}g^{bcy}g^{abcz}) , (a, b, c)]$$
(16)

The distribution (16) is not pseudorandom as any tuple  $(\alpha, \beta), (a, b, c) \in \mathcal{G}^2 \times \mathbb{Z}_p^3$  of the form (16) satisfies  $\alpha^c = \beta$ , which can be efficiently verified, while a random distribution will satisfy this relation only with probability 1/p.

#### 4.2 Expanding Seeds is Safe when Extracting Few Bits

By the following theorem, the expand-then-extract Hypothesis does hold, if the pseudorandom generator Prg used for expansion is sufficiently strong. The required hardness depends exponentially on the output length v of the extractor.

**Theorem 4.2** Assume Ext :  $\mathcal{X} \times \{0,1\}^n \to \{0,1\}^v$  is a  $(m,\varepsilon)$ -extractor with running time  $t_{\mathsf{Ext}}$ , and  $\mathsf{Prg} : \{0,1\}^k \to \{0,1\}^n$  is a  $(T,\sqrt{\varepsilon})$ -pseudorandom generator, for some

$$T \in O\left(2^{2v}(n+v)t_{\mathsf{Ext}}/\varepsilon\right) \tag{17}$$

Then  $\mathsf{Ext}'(x;s) \stackrel{\mathsf{def}}{=} \mathsf{Ext}(x; \mathsf{Prg}(s))$  is a  $(m, 4\sqrt{\varepsilon})$ -extractor. In particular, if the running time of  $\mathsf{Ext}$  is polynomial in k, its error  $\varepsilon(k) = \mathsf{negl}(k)$ , its output size  $v = O(\log k)$ , and  $\mathsf{Prg}$  is secure against polynomial (in k) size distinguishers, then  $\mathsf{Ext}'$  is an  $(m, \varepsilon')$ -extractor for  $\varepsilon'(k) = \mathsf{negl}(k)$ .

*Proof.* [(of Theorem 4.2)] Consider any distribution (X, Z) where  $\widetilde{\mathbf{H}}_{\infty}(X|Z) \geq m$ . Given fixed  $z \in \mathcal{Z}$ , let  $X_z$  be the distribution of X conditioned on Z = z. Also, for a given  $s' \in \{0, 1\}^n$  and  $z \in \mathcal{Z}$ , define

$$\phi(s', z) \stackrel{\text{def}}{=} \mathsf{SD}(\mathsf{Ext}(X_z; s') \ , \ U_v \ | Z = z \ , \ s' \ ) \tag{18}$$

We start with the following auxiliary claim which, roughly, states that the values  $\phi(s', z)$  can be approximated by circuits of size exponential in the number of extracted bits v:

**Claim 4.1** For any  $z \in \mathbb{Z}$  there is a circuit  $C_z$  of size  $T = O\left(2^{2v}(n+v)t_{\mathsf{Ext}}/\varepsilon\right)$  outputting a real number (up to some accuracy) such that, for all  $s' \in \{0,1\}^n$ , we have  $|C_z(s') - \phi(s',z)| < \frac{\sqrt{\varepsilon}}{2}$ .

*Proof.* [of Claim] Recall that

$$\phi(s',z) = \frac{1}{2} \sum_{w \in \{0,1\}^v} \left| \Pr[\mathsf{Ext}(X_z;s') = w] - \frac{1}{2^v} \right|$$
(19)

Thus, to prove the claim, it suffices to show the existence of circuit  $C'_z(s')$  of size O(T) which outputs  $2^v$  estimates  $\{p_w \mid w \in \{0,1\}^v\}$  for all probabilities  $\Pr[\mathsf{Ext}(X_z;s') = w]$ , where each  $p_w$  is within accuracy  $\delta \stackrel{\mathsf{def}}{=} \frac{\sqrt{\varepsilon}}{2} \cdot \frac{1}{2^v}$  from its true value:  $|p_w - \Pr[\mathsf{Ext}(X_z;s') = w]| < \delta$ .

Using the Chernoff bound, we observe that for any  $s' \in \{0,1\}^n$  and any  $w \in \{0,1\}^v$ , we can approximate the probability  $\Pr[\mathsf{Ext}(X_z;s') = w]$  with accuracy  $\delta$  and error probability strictly less than  $\tau \stackrel{\text{def}}{=} 2^{-n-v}$ , using "only"  $N \stackrel{\text{def}}{=} O(\log(1/\tau)/\delta^2) = O((n+v)2^{2v}/\varepsilon)$  independent samples of  $X_z$ . Taking the union bound over all such s' and w, we get that there exist N particular values of  $x_1, \ldots, x_N \in \mathcal{X}$  such that for all  $s' \in \{0,1\}^n$  and  $w \in \{0,1\}^v$ , all the empirical probabilities that  $\mathsf{Ext}(X_z;s') = w$  (computed by averaging the  $x_i$ 's above) are within  $\delta$  of their true values. Hardwiring these N values of  $x_1 \ldots x_N$  into our circuit  $C'_z(s')$  and running the extractor  $\mathsf{Ext}(x_i, s')$ on these values to compute all the empirical probabilities, we indeed get the required circuit  $C'_z(s')$ of size  $O(T) = O(Nt_{\mathsf{Ext}})$ .

Consider now a distinguisher  $D_z(s')$  defined as follows:

$$\mathcal{D}_z(s') = \begin{cases} 1 & \text{if } C_z(s') > \frac{3\sqrt{\varepsilon}}{2} \\ 0 & \text{otherwise} \end{cases}$$

Since  $\mathcal{D}_z$  has size T and our PRG is  $(T, \sqrt{\varepsilon})$ -secure, we have

$$\Pr[\mathcal{D}_z(\mathsf{Prg}(U_k)) = 1] \le \Pr[\mathcal{D}_z(U_n) = 1] + \sqrt{\varepsilon}$$

Recalling the definition of  $\mathcal{D}_z$  above, and that the circuit  $C_z$  approximates  $\phi(s', z)$  with accuracy  $\sqrt{\varepsilon}/2$  (by Claim), we get that for any value of  $z \in \mathbb{Z}$ ,

$$\Pr[\phi(\Pr(U_k), z) > 2\sqrt{\varepsilon}] \leq \Pr\left[C_z(\Pr(U_k)) > \frac{3\sqrt{\varepsilon}}{2}\right] = \Pr[\mathcal{D}_z(\Pr(U_k)) = 1]$$
$$\leq \Pr[\mathcal{D}_z(U_n) = 1] + \sqrt{\varepsilon} = \Pr\left[C_z(U_n) > \frac{3\sqrt{\varepsilon}}{2}\right] + \sqrt{\varepsilon}$$
$$\leq \Pr[\phi(U_n, z) > \sqrt{\varepsilon}] + \sqrt{\varepsilon}$$

Taking the expectation over Z, we get

$$\Pr[\phi(\mathsf{Prg}(U_k), Z) > 2\sqrt{\varepsilon}] \le \Pr[\phi(U_n, Z) > \sqrt{\varepsilon}] + \sqrt{\varepsilon}$$
(20)

Recalling now the definition of  $\phi(s', z)$  from Equation (18), and using Equation (20), we get

$$\begin{aligned} \mathsf{SD}(\mathsf{Ext}'(X;S) \ , \ U_v \ | Z \ , \ S \ ) &= & \mathsf{SD}(\mathsf{Ext}(X;\mathsf{Prg}(S)) \ , \ U_v \ | Z \ , \ S \ ) \\ &= & \mathbb{E}_{s \leftarrow U_k, z \leftarrow Z}[\phi(\mathsf{Prg}(s), z)] \\ &\leq & \operatorname{Pr}[\phi(\mathsf{Prg}(U_k), Z) > 2\sqrt{\varepsilon}] \cdot 1 + 1 \cdot 2\sqrt{\varepsilon} \\ &\leq & \operatorname{Pr}[\phi(U_n, Z) > \sqrt{\varepsilon}] + 3\sqrt{\varepsilon} \end{aligned}$$

Finally, because Ext is  $(m, \varepsilon)$ -extractor, we must have

$$\varepsilon \geq \mathsf{SD}(\mathsf{Ext}(X;S') \ , \ U_v \ | Z \ , \ S' \ ) = \mathbb{E}_{s',z}[\phi(s',z)] \geq \Pr[\phi(U_n,Z) > \sqrt{\varepsilon}] \cdot \sqrt{\varepsilon}$$
  
which means that  $\Pr[\phi(U_n,Z) > \sqrt{\varepsilon}] \leq \sqrt{\varepsilon}$  and  $\mathsf{SD}(\mathsf{Ext}'(X;S) \ , \ U_v \ | Z \ , \ S \ ) \leq 4\sqrt{\varepsilon}.$ 

DISCUSSION. It is interesting to examine the best possible seed length one can obtain by using Theorem 4.2, even with exponentially-secure PRGs. Replacing  $\varepsilon$  by  $\sqrt{\varepsilon}$ , we get that in order to obtain  $(m, \varepsilon)$ -secure v-bit extractor, we need a PRG which is  $(T, O(\varepsilon))$ -secure against circuits of size  $\tilde{O}(4^v/\varepsilon^2)$ . Clearly, such a PRG must have a seed of length  $k \ge \log T = \Omega(v + \log(1/\varepsilon))$ . Conversely, under an exponential assumption on the PRG, such a seed length suffices. Assuming a PRG with such a seed length is not unreasonable, provided that the length of the extracted key v is not too large. I.e., for v = 128 and  $\varepsilon = 2^{-80}$ , it seems plausible that a seed length of size, say, 1000 might already suffice. The bad news is that we can achieve the same seed length  $O(v + \log(1/\varepsilon))$ using *almost* universal hash functions [40, 39], without any computational assumptions, and without too much (if any) efficiency degradation. Thus, the result of Theorem 4.2 is mainly of theoretical interest, since better constructions for the same goal are easily achievable.

#### 4.3 Expanding Seeds is Safe in Minicrypt

Before we can state the main result of this section, we need a few more definitions.

BIT-AGREEMENT. Bit-agreement is a protocol between two efficient parties, which we refer to as Alice and Bob. They get the security parameter k in unary (denoted  $1^k$ ) as a common input and can communicate over an authentic channel. Finally, Alice and Bob output a bit  $b_A$  and  $b_B$ , respectively. The protocol has correlation  $\epsilon = \epsilon(k)$ , if for all k,  $\Pr[b_A = b_B] \ge (1 + \epsilon(k))/2$ . Furthermore, the protocol has security  $\delta = \delta(k)$ , if for every efficient adversary Eve, which can observe the whole communication C, and for all k,  $\Pr[\operatorname{Eve}(1^k, C) = b_B] \le 1 - \delta(k)/2$ .

KEY-AGREEMENT & PKE. If  $\epsilon(\cdot)$  and  $\delta(\cdot)$  are overwhelming then such a protocol achieves *key*agreement. Using parallel repetition and privacy amplification, it is known [21, 20] that any protocol which achieves bit-agreement with noticeable correlation  $\epsilon(\cdot)$  and overwhelming security  $\delta(\cdot)$  can be turned into a key-agreement protocol, without increasing the number of rounds. A 2-message key-agreement protocol is equivalent to public-key encryption (PKE).

**Theorem 4.3 (Hypothesis 1 holds in minicrypt)** If there exists a secure pseudorandom generator Prg and a strong extractor Ext where  $Ext'(.;.) \stackrel{\text{def}}{=} Ext(Prg(.);.)$  is <u>not</u> a computational extractor, then the protocol from Figure 1 is a two-message bit-agreement protocol with noticeable correlation and overwhelming security (and thus implies PKE).

**Remark 4.1** In the above theorem not being a secure computational extractor means that there exists an efficient uniform D that can distinguish Ext(Prg(.);.) with <u>noticeable</u> advantage (in the security parameter k). If Ext(Prg(.);.) is only insecure against non-uniform adversaries, then also

the resulting protocol (which uses D) will be non-uniform. If the distinguisher D only has nonnegligible advantage (i.e. only works for infinitely many, but not all, security parameters k), then also the protocol will work for infinitely many k. This issue is inherent in win-win type results where an adversary is turned into a "useful" protocol [13, 34, 11, 35]. It roots in the fact that in cryptography we usually put weak requirements for adversaries to be considered efficient (can be non-uniform and only have non-negligible advantage), whereas we usually require from practical algorithms to be uniform and secure for all (sufficiently large) security parameters.

*Proof.* [of Theorem 4.3] Let  $\mathsf{Ext} : \mathcal{X} \times \{0,1\}^n \to \{0,1\}^v$ ,  $\mathsf{Prg} : \{0,1\}^k \to \{0,1\}^n$  be as in the Theorem, i.e. where for some efficiently samplable distribution (X,Z) s.t.  $\widetilde{\mathbf{H}}_{\infty}(X|Z) \ge m = m(k)$ , there exists a distinguisher D such that for some noticeable  $\mu = \mu(k)$ , with  $s \leftarrow U_k$ ,  $r \leftarrow U_v$ 

$$\Pr[D(\mathsf{Ext}(x; \mathsf{Prg}(s)), s, z) = 1] - \Pr[D(r, s, z) = 1] \ge \mu(k)$$
(21)

#### Claim 4.2 The protocol has overwhelming security.

*Proof.* [of Claim] An adversary Eve observing the communication between Alice and Bob sees either

$$Z_0 \equiv (\mathsf{Ext}(x; \mathsf{Prg}(s)), \mathsf{Prg}(s), z) \quad \text{if} \quad b_{\mathsf{B}} = 0 \quad \text{or} \quad Z_1 \equiv (U_v, \mathsf{Prg}(s), z) \quad \text{if} \quad b_{\mathsf{B}} = 1$$

It follows directly from the security of Prg and Ext that  $Z_0$  and  $Z_1$  are pseudorandom, and thus computationally indistinguishable. I.e.

$$\Pr[\mathsf{Eve}(1^k, \{r, s', z\}) = b_{\mathsf{B}}] = 1/2 + \mathsf{negl}(k) = 1 - \frac{\delta(k)}{2}$$

for an overwhelming  $\delta(k)$ .

**Claim 4.3** The protocol has noticeable correlation  $\mu(k)$ .

*Proof.* [of Claim] By Equation (21) and the definition of the protocol (where  $b_{\mathsf{A}} = D(s, r)$  with  $r \equiv \mathsf{Ext}(x; \mathsf{Prg}(s))$  if  $b_{\mathsf{B}} = 0$  and  $r \equiv U_v$  if  $b_{\mathsf{B}} = 1$ ) we get

$$\Pr[b_{\mathsf{A}} = 1 | b_{\mathsf{B}} = 1] - \Pr[b_{\mathsf{A}} = 1 | b_{\mathsf{B}} = 0] \ge \mu$$

Which implies noticeable correlation  $\mu$ 

$$\Pr[b_{\mathsf{A}} = b_{\mathsf{B}}] = \underbrace{\Pr[b_{\mathsf{B}} = 1]}_{=1/2} \Pr[b_{\mathsf{A}} = 1 | b_{\mathsf{B}} = 1] + \underbrace{\Pr[b_{\mathsf{B}} = 0]}_{=1/2} \underbrace{\Pr[b_{\mathsf{A}} = 0 | b_{\mathsf{B}} = 0]}_{=1 - \Pr[b_{\mathsf{A}} = 1 | b_{\mathsf{B}} = 0]}$$

$$\geq \frac{1}{2} \left(\Pr[b_{\mathsf{A}} = 1 | b_{\mathsf{B}} = 1] + 1 - \Pr[b_{\mathsf{A}} = 1 | b_{\mathsf{B}} = 0]\right) = \frac{1 + \mu}{2}$$

We obtain the following corollary whose proof is immediate from Figure 1 and our proof of Theorem 4.3.

$$\begin{array}{cccc} \text{Alice} & \text{Bob} \\ s \leftarrow U_k; s' \leftarrow \mathsf{Prg}(s) & \xrightarrow{s'} & (x, z) \leftarrow (X, Z) \\ & & & b_{\mathsf{B}} \leftarrow U_1 \\ & & & \text{if } b_{\mathsf{B}} = 0 \text{ then } r \leftarrow \mathsf{Ext}(x; s') \\ & & & \text{else if } b_{\mathsf{B}} = 1 \text{ then } r \leftarrow U_v \end{array}$$

Figure 1: A bit agreement protocol from any Prg, Ext that constitute a counterexample (via distinguisher D) to Hypothesis 1.

**Corollary 4.1** Assume Prg is a secure pseudorandom generator. Assume further that there exists no public-key encryption scheme (with non-negligible gap between security and decryption correctness) having: (a) pseudorandom ciphertexts, (b) secret key equal to the seed of Prg and public key equal to the output of Prg; (c) tight security reduction to the security of Prg. Then the expand-then-extract hypothesis is true for Prg.

DISCUSSION. Corollary 4.1 reduces the soundness of the expand-then-extract approach to the impossibility of constructing public-key encryption which is as secure as the given PRG, and, moreover, has a very particular form: its ciphertexts are pseudorandom and, more importantly, the key-generation algorithm samples a random s and sets sk = s, pk = Prg(s). Of course, specific number-theoretic PRGs, such as the PRG from Equation (15) used in our counter-example, might very well give such an encryption scheme. E.g., a simple variant of ElGamal encryption scheme would be the desired encryption scheme for the DDH-based PRG used in our counter-example. However, to the best of our knowledge, this impossibility assumption seems very likely for "practical" PRGs, such as AES (in counter mode, say), which are not based on number theory and do not seem to have any algebraic structure. For example, not only there is no black-box construction of PKE (or key agreement) from a PRG alone, as shown by Impagliazzo and Rudich [24], but, in fact, it is entirely consistent with current knowledge that these two tasks are separable, in the sense that there is some computational model/complexity class (e.g., perhaps some extension of BQP or  $AM \cap coAM$ ) that is powerful enough to break all public key schemes, but not powerful enough to break AES. If this is the case, then the AES-bases expand and extract scheme is secure with respect to all efficient input distributions and distinguishers (even those that are based on public key tools such as factoring, lattices etc.), since the resulting (hypothetical) PKE cannot be as secure as AES. Moreover, we do not know any black-box construction of PKE from PRG and any other "cryptomania" assumption (like non-interactive zero-knowledge proofs, fully-homomorphic or identity-based encryption, etc.), where the public key of the PKE is simply the output of the PRG, even if the reduction is allowed to be loose (which is not allowed in our case). To summarize, our results give strong evidence that the expand-then-extract approach is secure using any "practical" PRG (like AES), despite being generally insecure in theory (e.g., when used with "number-theoretic" PRGs).

CONCRETE SECURITY. Assuming one believes the assumptions of Corollary 4.1, we can try to informally estimate the concrete (computational) security  $\varepsilon'$  of Ext' as a function of (statistical) security  $\varepsilon$  of Ext and (computational) security  $\delta$  of Prg. Examining the proof of Theorem 4.3, we can see that the key agreement protocol achieves (computational) semantic security  $\varepsilon + \delta$ , and achieves a non-negligible gap between this security and correctness of decryption as long as the hypothetical distinguisher D breaks the expand-then-extract approach with advantage slightly more than  $\varepsilon + \delta$ . Thus, it is reasonable to conjecture that, with a "friendly" PRG, Ext' achieves concrete security only negligibly worse than  $\varepsilon + \delta$ , which was precisely the naive hope one might have had at the beginning.

OBLIVIOUS TRANSFER. In fact, notice that the PKE that we obtained from our PRG has the following property: if one replaces the public key of the PKE from pseudorandom  $pk = \Pr(s)$  to truly random  $pk \leftarrow U_n$ , then the encrypted bit b is statistically hidden. It is easy to see that such a special PKE easily yields a 2-round *oblivious transfer* protocol secure against honest-but-curious receiver, which is one of the the strongest primitives in cryptomania. Indeed, imagine Alice has a choice bit  $b \in \{0, 1\}$ . She then sets  $pk_b \leftarrow \Pr(S)$  and  $pk_{1-b} \leftarrow U_n$  and sends  $pk_0, pk_1$  to Bob. Bob, who holds two bits  $\sigma_0$  and  $\sigma_1$ , uses our PKE to encrypt  $\sigma_0$  and  $\sigma_1$  under  $pk_0$  and  $pk_1$ , respectively. Alice can then recover  $\sigma_b$  and gets no information about  $\sigma_{1-b}$ .

**Remark 4.2** So far we stated the expand-then-extract approach assuming that the original "longseed" extractor Ext is a statistical extractor. However, everything we said also holds if Ext itself is a "long-seed" computational extractor, such as the one designed in Section 3.3.

### 5 Extractor Standard

Although we leave the detailed exploration of this direction to future work, one can easily combine the results of Section 3 and Section 4 to get a concrete randomness extractor — leading perhaps to an "extractor standard" advocated in the Introduction — with an improved entropy loss and a fixed-length seed. The idea is to use a pseudorandom function F, such as AES, to first implement the general-purpose (long-seed) computational extractor from Section 3.3, and then use the expandthen-extract approach, with the same PRF in the counter mode, to provide the long seed required by this computational extractor. As explained in Remark 4.2, this approach is still sound and will provide the improved entropy loss bounds with a fixed-length seed, as long as one believes the assumptions stated in Corollary 4.1. Moreover, under proper optimizations, the efficiency of our approach will likely be comparable with (or faster than) that of existing cryptographic hash functions.

Of course, figuring out the optimal (almost) universal hash function, as well as the best setting of parameters to compose such a hash functions with a practical PRF, requires non-trivial engineering effort, and will likely depend on a particular scenario as hand (e.g., hardware support for fast field arithmetic, availability of cryptographic co-processor, etc.) Therefore, we leave this interesting research to future work, here only providing an *extremely simple* implementation of such an extractor, based on the AES block cipher and the simple inner-product universal hash function. Although our implementation might not be the fastest way to instantiate our general approach, its efficiency is already compatible with the speed of existing cryptographic hash functions.

CONCRETE PROPOSAL. We will use the AES block cipher which has key, input and output size equal to k = 128 bits, and use the notation  $AES_s(w)$  to denote the evaluation of AES with key s and input w. Given a (variable-length) n-block input  $X = (x_1 \dots x_n)$ , we let  $x_{n+1} = \langle n \rangle$ denote a block holding the length of X, which we append to X to ensure suffix-freeness. We will also use the following simple inner product function IP(X;S), whose (variable-length) input  $X = (x_1 \dots x_n, x_{n+1})$  and seed  $S = (s_1 \dots s_{n+1})$  are both interpreted as n elements of the finite field  $GF[2^{128}]$  (thus, addition '+' and multiplication '.' are interpreted according to the corresponding finite field operation):

$$\mathsf{IP}(X;S) = \mathsf{IP}((x_1, x_2, \dots, x_n, x_{n+1}); (s_1, s_2, \dots, s_{n+1})) \stackrel{\text{def}}{=} x_1 \cdot s_1 + x_2 \cdot s_2 + \dots + x_n \cdot s_n + x_{n+1} \cdot s_{n+1}$$
(22)

It is well known that the above inner-product function is perfectly universal (so  $\gamma = 0$ ) even for variable-length inputs X (this is why we set  $x_{n+1} = \langle n \rangle$ ). We will first apply the technique from

Section 3.3, and use a seed S' = (w, S), where  $w \in GF[2^{128}]$  is a random point required for the wPRF application. Namely, we define

$$\mathsf{Ext}'(X;S') = \mathsf{Ext}'(X;(w,S)) \stackrel{\text{def}}{=} \mathsf{AES}_{\mathsf{IP}(X;S)}(w) = \mathsf{AES}_{x_1 \cdot s_1 + \dots + x_{n+1} \cdot x_{n+1}}(w)$$
(23)

Finally, we will use the expand-then-extract approach to expand a single 128-bit seed s to generate the required (n + 2) blocks of the long seed S', and define our final computational extractor as follows:

$$\mathsf{Ext}(X;s) \stackrel{\text{def}}{=} \mathsf{Ext}'(X;(\mathsf{AES}_s(0),\mathsf{AES}_s(1),\dots,\mathsf{AES}_s(n)))$$
(24)

$$= \mathsf{AES}_{x_1 \cdot \mathsf{AES}_s(1) + \dots + x_{n+1} \cdot \mathsf{AES}_s(n+1)}(\mathsf{AES}_s(0))$$
(25)

SUPPORTING LONGER OUTPUTS. We notice that our final extractor Ext only extracts k = 128 bits. Of course, if the application needs to extract v > 128 bits, we can always use a pseudorandom generator to expand our extracted 128-bit key to produce the desired number of bits. E.g., in the case of using AES, we can simply use the value r = Ext(X;s) and output  $\text{AES}_r(1), \ldots, \text{AES}_r(t)$ , where  $t = \lceil v/128 \rceil$ . Assuming the counter mode of AES is a good stream cipher (which is anyway needed to apply the expand-then-extract approach), the resulting v bits will be pseudorandom as long as the initial key r is such.

EFFICIENCY NOTES. We notice that the key schedule for AES can be reused and (pre-)computed only once when evaluating  $AES_s(0)$ ,  $AES_s(1)$ ,.... Similar reuse of the key schedule can also be done when the output of our extractor needs to be longer than 128 bits, as explained above. Finally, our computational extractor is friendly to variable-length streaming sources  $X = x_1, x_2, \ldots$ , as the intermediate AES key  $x_1 \cdot AES_s(1) + \ldots + x_{n+1} \cdot AES_s(n+1)$  can be easily computed on-line, without a-priori knowing the number of blocks n.

CONCRETE SECURITY. Assume the min-entropy of X is m = 128 + L, where L is the entropy loss. Consider also some application P which requires a random v-bit key r, and has ideal security  $\delta$  in this case (against some class of attackers). For simplicity, we assume v = 128, as longer outputs can be computed pseudorandomly from a 128-bit output, as explained above. We now derive bounds of the "real" security  $\delta'$  of P when we use the extracted values IP(X; S) (see Equation (22)), Ext'(X; S') (see Equation (23)) and Ext(X; s) (see Equation (25)), respectively, as the "real" key for P.

Using  $\mathsf{IP}(X; S)$ . Applying the standard LHL to the perfectly universal inner product function  $\mathsf{IP}$ , we get  $\delta' \leq \delta + \sqrt{2^{-L}}$ . For the common case when P is a "GLHL-friendly" application (e.g., MAC or stateless CPA encryption scheme), the results form Section 3.2 give an improved bound

$$\delta' \le \delta + \sqrt{\delta \cdot 2^{-L}} \tag{26}$$

Using Ext'(X; S'). Here we no longer need to assume that P is "GLHL"-friendly, and can use the computational results from Section 3.3. Namely, let  $\varepsilon$  be the wPRF security of AES (against the same complexity attackers as P) against a single wPRF query, as defined in Theorem 3.7. Then, Theorem 3.7 implies that

$$\delta' \le \delta + \varepsilon + \sqrt{\varepsilon 2^{-L}} + \frac{1}{2^{128}} \tag{27}$$

This bound has two advantages over the bound in Equation (26), when using IP(X; S) as the key. First, we no longer need to settle for "GLHL"-friendly applications (e.g., P could be a

stream cipher). Second, in practice we often expect  $\delta \gg \varepsilon$ , since  $\varepsilon$  is a the security of AES against very limited wPRF attackers (who only learn the real AES value at one random point). In this case, to achieve  $\delta' = O(\delta)$ , the previous analysis requires entropy loss  $L = \log(1/\delta)$ , while the new analysis only needs entropy loss  $L = 2\log(1/\delta) - \log(1/\varepsilon) < \log(1/\delta)$ , as  $\delta \gg \varepsilon$ .

Using  $\operatorname{Ext}(X; s)$ . The previous two extractors IP and  $\operatorname{Ext}'$  had long (potentially variable-length) seeds S and S' = (S, w). Our final extractor  $\operatorname{Ext}$  uses a *short* 128-*bit seed* s, and will be a good computational extractor assuming AES satisfies the conditions of Corollary 4.1.<sup>10</sup> Specifically, let  $\varepsilon'$  be the security of AES as a stream cipher in the counter mode (against the same complexity attackers as P). In practice we suspect that  $\varepsilon' \approx \varepsilon$  (where, recall,  $\varepsilon$  is the wPRF security of AES). As argued at the end of Section 4.3, although it is hard to precisely assess the security of Ext, it seems reasonable to assume that its computational security  $\delta'$  is only "negligibly worse" than  $\varepsilon'$  plus the computational security of Ext' given in Equation (27). In particular, it seems reasonable to conjecture that

$$\delta' = O(\delta + \varepsilon' + \varepsilon + \sqrt{\varepsilon 2^{-L}}) \tag{28}$$

Since the above bound already includes the AES stream cipher security  $\varepsilon'$ , the same bound also holds when the extracted key  $\mathsf{Ext}(X;s)$  is used in the stream cipher mode to output more than 128 output bits.

Acknowledgements: We would like to thank Russell Impagliazzo, Ronen Shaltiel and Daniel Wichs for useful discussions, and Gil Segev for pointing out the our construction in Section 4.3 also implies oblivious transfer.

# References

- Boaz Barak and Shai Halevi. A model and architecture for pseudo-random generation with applications to /dev/random. In Vijay Atluri, Catherine Meadows, and Ari Juels, editors, ACM Conference on Computer and Communications Security, pages 203–212. ACM, 2005.
- [2] Boaz Barak, Ronen Shaltiel, and Eran Tromer. True random number generators secure in a changing environment. In *Cryptographic Hardware and Embedded Systems – CHES '03*, volume 2779 of *LNCS*, pages 166–180. Springer, 2003.
- [3] Charles H. Bennett, Gilles Brassard, and Jean-Marc Robert. Privacy amplification by public discussion. SIAM Journal on Computing, 17(2):210–229, 1988.
- [4] Xavier Boyen, Yevgeniy Dodis, Jonathan Katz, Rafail Ostrovsky, and Adam Smith. Secure remote authentication using biometric data. In Ronald Cramer, editor, Advances in Cryptology— EUROCRYPT 2005, volume 3494 of LNCS, pages 147–163. Springer-Verlag, 2005.
- [5] Ran Canetti, Yevgeniy Dodis, Shai Halevi, Eyal Kushilevitz, and Amit Sahai. Exposureresilient functions and all-or-nothing transforms. In Bart Preneel, editor, Advances in Cryptology—EUROCRYPT 2000, volume 1807 of LNCS, pages 453–469. Springer-Verlag, 2000.
- [6] J.L. Carter and M.N. Wegman. Universal classes of hash functions. Journal of Computer and System Sciences, 18:143–154, 1979.

<sup>&</sup>lt;sup>10</sup>As explained in Remark 4.2, Corollary 4.1 also applies when the original long-seed extractor is computationally secure, such as Ext' above.

- [7] Céline Chevalier, Pierre-Alain Fouque, David Pointcheval, and Sébastien Zimmer. Optimal randomness extraction from a diffie-hellman element. In Antoine Joux, editor, Advances in Cryptology - EUROCRYPT 2009, volume 5479 of LNCS, pages 572–589. Springer-Verlag, 2009.
- [8] Yevgeniy Dodis, Rosario Gennaro, Johan Håstad, Hugo Krawczyk, and Tal Rabin. Randomness extraction and key derivation using the cbc, cascade and hmac modes. In Franklin [14], pages 494–510.
- [9] Yevgeniy Dodis, Jonathan Katz, Leonid Reyzin, and Adam Smith. Robust fuzzy extractors and authenticated key agreement from close secrets. In Dwork [12], pages 232–250.
- [10] Yevgeniy Dodis, Rafail Ostrovsky, Leonid Reyzin, and Adam Smith. Fuzzy extractors: How to generate strong keys from biometrics and other noisy data. SIAM Journal on Computing, 38(1):97–139, 2008.
- [11] Bella Dubrov and Yuval Ishai. On the randomness complexity of efficient sampling. In STOC, pages 711–720, 2006.
- [12] Cynthia Dwork, editor. Advances in Cryptology—CRYPTO 2006, volume 4117 of LNCS. Springer-Verlag, 20–24 August 2006.
- [13] Stefan Dziembowski. On forward-secure storage. In Dwork [12], pages 251–270.
- [14] Matt Franklin, editor. Advances in Cryptology—CRYPTO 2004, volume 3152 of LNCS. Springer-Verlag, 15–19 August 2004.
- [15] Rosario Gennaro, Hugo Krawczyk, and Tal Rabin. Secure hashed diffie-hellman over nonddh groups. In Christian Cachin and Jan Camenisch, editors, Advances in Cryptology— EUROCRYPT 2004, volume 3027 of LNCS, pages 361–381. Springer-Verlag, 2004.
- [16] O. Goldreich and L. Levin. A hard-core predicate for all one-way functions. In Johnson [26], pages 25–32.
- [17] Venkatesan Guruswami, Christopher Umans, and Salil P. Vadhan. Unbalanced expanders and randomness extractors from parvaresh-vardy codes. J. ACM, 56(4), 2009.
- [18] Gustav Hast. Nearly one-sided tests and the goldreich?levin predicate. J. Cryptology, 17(3):209–229, 2004.
- [19] J. Håstad, R. Impagliazzo, L.A. Levin, and M. Luby. Construction of pseudorandom generator from any one-way function. SIAM Journal on Computing, 28(4):1364–1396, 1999.
- [20] Thomas Holenstein. Key agreement from weak bit agreement. In STOC, pages 664–673, Baltimore, Maryland, 22–24 May 2005.
- [21] Thomas Holenstein. Strengthening Key Agreement using Hard-Core Sets. PhD thesis, ETH Zurich, Zurich, Switzerland, 2006.
- [22] Chun-Yuan Hsiao and Leonid Reyzin. Finding collisions on a public road, or do secure hash functions need secret coins. In Franklin [14], pages 92–105.
- [23] Russell Impagliazzo. A personal view of average-case complexity. In Structure in Complexity Theory Conference, pages 134–147, 1995.

- [24] Russell Impagliazzo and Steven Rudich. Limits on the provable consequences of one-way permutations. In Johnson [26], pages 44–61.
- [25] Yuval Ishai, Eyal Kushilevitz, Rafail Ostrovsky, and Amit Sahai. Cryptography with constant computational overhead. In Richard E. Ladner and Cynthia Dwork, editors, *Proceedings of* the 40th Annual ACM Symposium on Theory of Computing, pages 433–442, Victoria, BC, Canada, 17–20 May 2008. ACM.
- [26] D. S. Johnson, editor. Proceedings of the Twenty First Annual ACM Symposium on Theory of Computing, Seattle, Washington, 15–17 May 1989.
- [27] Hugo Krawczyk. Cryptographic Extraction and Key Derivation: The HKDF Scheme. In Tal Rabin, editor, Advances in Cryptology - CRYPTO 2010, volume 6223 of LNCS, pages 631–648. Springer-Verlag, 2010.
- [28] Ueli Maurer and Stefan Wolf. Privacy amplification secure against active adversaries. In Burton S. Kaliski, Jr., editor, Advances in Cryptology—CRYPTO '97, volume 1294 of LNCS, pages 307–321. Springer-Verlag, 1997.
- [29] Michael Mitzenmacher and Salil Vadhan. Why simple hash functions work: Exploiting the entropy in a data stream. In In Proceedings of the 19th Annual ACM-SIAM Symposium on Discrete Algorithms, pages 746–755, 2008.
- [30] Moni Naor and Omer Reingold. Number-theoretic constructions of efficient pseudo-random functions. In 38th Annual Symposium on Foundations of Computer Science, pages 458–467, Miami Beach, Florida, 20–22 October 1997. IEEE.
- [31] Moni Naor and Gil Segev. Public-key cryptosystems resilient to key leakage. In Shai Halevi, editor, Advances in Cryptology - CRYPTO 2009, volume 5677 of LNCS, pages 18–35. Springer-Verlag, 2009.
- [32] Wim Nevelsteen and Bart Preneel. Software performance of universal hash functions. In Jacques Stern, editor, Advances in Cryptology—EUROCRYPT '99, volume 1592 of LNCS, pages 24–41. Springer-Verlag, 2–6 May 1999.
- [33] Noam Nisan and David Zuckerman. Randomness is linear in space. Journal of Computer and System Sciences, 52(1):43–53, 1996.
- [34] Krzysztof Pietrzak. Composition implies adaptive security in minicrypt. In Serge Vaudenay, editor, *EUROCRYPT*, volume 4004 of *LNCS*, pages 328–338. Springer-Verlag, 2006.
- [35] Krzysztof Pietrzak and Johan Sjödin. Weak pseudorandom functions in minicrypt. In Luca Aceto, Ivan Damgård, Leslie Ann Goldberg, Magnús M. Halldórsson, Anna Ingólfsdóttir, and Igor Walukiewicz, editors, *ICALP (2)*, volume 5126 of *Lecture Notes in Computer Science*, pages 423–436. Springer, July 7-11 2008.
- [36] Jaikumar Radhakrishnan and Amnon Ta-Shma. Bounds for dispersers, extractors, and depthtwo superconcentrators. SIAM Journal on Computing, 13(1):2–24, 2000.
- [37] Renato Renner and Stefan Wolf. Unconditional authenticity and privacy from an arbitrarily weak secret. In Dan Boneh, editor, Advances in Cryptology—CRYPTO 2003, volume 2729 of LNCS, pages 78–95. Springer-Verlag, 2003.

- [38] Ronen Shaltiel. Recent developments in explicit constructions of extractors. Bulletin of the EATCS, 77:67–95, 2002.
- [39] A. Srinivasan and D. Zuckerman. Computing with very weak random sources. SIAM J. Computing, 28(4):1433–1459, 1999.
- [40] D. R. Stinson. Universal hashing and authentication codes. Designs, Codes, and Cryptography, 4(4):369–380, 1994.
- [41] D. Stinson. families R. Universal hashand the leftover hash lemma, and applications to cryptography and computing. Journal ofCombinatorial Mathematics and Combinatorial Computing, 42:3-31,2002. Available at http://www.cacr.math.uwaterloo.ca/~dstinson/publist.html.

# A Slice Extractors

For completeness, we also show that universal hash function also give nearly optimal *slice extractors* — a primitive defined by [36]. This application is similar to our improvements of the LHL for authentication application, but we use a slightly different notation (based on that of [36]).

We say that a distinguisher D is  $\delta$ -rare on a distribution X, if  $\Pr[D(X) = 1] \leq \delta$ . We recall the definition of *slice extractors* from [36].

**Definition A.1 (Slice Extractors)** We say that an efficient function  $\mathsf{Ext} : \mathcal{X} \times \{0,1\}^n \to \{0,1\}^v$ is an (average-case, strong)  $(m,\varepsilon,\delta)$ -slice extractor (for space  $\mathcal{X}$ ), if for all X, Z such that X is distributed over  $\mathcal{X}$  and  $\widetilde{\mathbf{H}}_{\infty}(X|Z) \geq m$ , and all  $\delta$ -rare distinguishers D on  $(U_v, S, Z)$ , we have

$$\Delta_D(\operatorname{Ext}(X;S), U_v \mid (S,Z)) \le \varepsilon$$

where  $S \equiv U_n$  denotes the coins of Ext (called the seed). In particular,

$$\Pr[D(\mathsf{Ext}(X;S),S,Z)=1] \le \Pr[D(U_v,S,Z)=1] + \varepsilon \le \delta + \varepsilon$$

As usual, the value L = m - v is called the entropy loss of Ext, and the value n is called the seed length of Ext.

Radhakrishnan and Ta-Shma showed (Lemma 2.5 in [36]) that the entropy loss of  $(m, \varepsilon, \delta)$ -slice extractors must be at least  $L \geq 2 \log (1/\varepsilon) - \log (1/\delta) - O(1)$ . We now state our improved variant of LHL for slice extractors, which constructively matches the bound above.

**Lemma A.1 ("Sliced LHL")** Assume that the family  $\mathcal{H}$  of functions  $h : \mathcal{X} \to \{0,1\}^v$  is a  $\rho$ -universal hash family, where  $\rho \leq \frac{1}{2^v}(1 + \varepsilon^2/\delta)$ . Then the extractor  $\mathsf{Ext}(x;h) \stackrel{\text{def}}{=} h(x)$ , where h is uniform over  $\mathcal{H}$ , is an  $(m, \varepsilon, \delta)$ -slice extractor, as long as  $m \geq v + 2\log(1/\varepsilon) - \log(1/\delta) + 1$ .

Hence,  $\rho$ -universal hash functions yield slice extractors with entropy loss  $L = 2\log(1/\varepsilon) - \log(1/\delta) + 1$  and seed length  $n = \log |\mathcal{H}|$ . In particular, when  $\delta = \varepsilon$ , we get  $(m, \varepsilon, \varepsilon)$ -slice extractor with entropy loss  $L = \log(1/\varepsilon) + 1$ .

*Proof.* Follows from part (b) of Theorem 3.2 by renaming variables.