

Pseudorandom Correlation Breakers, Independence Preserving Mergers and their Applications

Xin Li *

Department of Computer Science
 Johns Hopkins University
 Baltimore, MD 21218, U.S.A.
 lixints@cs.jhu.edu

Abstract

The recent line of study on randomness extractors has been a great success, resulting in exciting new techniques, new connections, and breakthroughs to long standing open problems in the following five seemingly different topics: seeded non-malleable extractors, privacy amplification protocols with an active adversary, independent source extractors (and explicit Ramsey graphs), non-malleable independent source extractors, and non-malleable codes in the split state model. Two key ingredients used in these works are correlation breakers and independence preserving mergers. By giving very efficient constructions of these two objects, we now have close to optimal solutions to the above five problems [Li17]: seeded non-malleable extractors with seed length and entropy requirement $O(\log n + \log(1/\epsilon) \log \log(1/\epsilon))$ for error ϵ ; two-round privacy amplification protocols with optimal entropy loss for security parameter up to $\Omega(k/\log k)$, where k is the entropy of the shared weak source; two-source extractors for entropy $O(\log n \log \log n)$; non-malleable two-source extractors for entropy $(1 - \gamma)n$ with error $2^{-\Omega(n/\log n)}$; and non-malleable codes in the 2-split state model with rate $\Omega(1/\log n)$. However, in all cases there is still a small gap to optimum and the motivation to close this gap remains strong. On the other hand, previous techniques seem to have reached their limit and insufficient for this purpose.

In this paper we introduce new techniques to recycle the entropy used in correlation breakers and independence preserving mergers. This allows us to break the barriers of previous techniques and give further improvements to the above problems. Specifically, we obtain the following results: (1) a seeded non-malleable extractor with seed length $O(\log n) + \log^{1+o(1)}(1/\epsilon)$ and entropy requirement $O(\log \log n + \log(1/\epsilon))$, where the entropy requirement is asymptotically optimal by a recent result of Gur and Shinkar [GS18]; (2) a two-round privacy amplification protocol with optimal entropy loss for security parameter up to $\Omega(k)$, which solves the privacy amplification problem completely;¹ (3) a two-source extractor for entropy $O(\frac{\log n \log \log n}{\log \log \log n})$, which also gives an explicit Ramsey graph on N vertices with no clique or independent set of size $(\log N)^{O(\frac{\log \log \log N}{\log \log \log \log N})}$; (4) a non-malleable two-source extractor for entropy $(1 - \gamma)n$ with error $2^{-\Omega(n \log \log n / \log n)}$; and (5) non-malleable codes in the 2-split state model with rate $\Omega(\log \log n / \log n)$. Some of our techniques are similar in spirit to what has been done in previous constructions of pseudorandom generators for small space computation [Nis92, NZ96], and we believe they can be a promising way to eventually obtain optimal constructions to the five problems mentioned above.

*Supported by NSF award CCF-1617713.

¹Except for the communication complexity, which is of secondary concern to this problem.

1 Introduction

The study of randomness extractors has been a central line of research in the area of pseudorandomness, where the goal is to understand how to use randomness more efficiently in computation. As fundamental objects in this area, randomness extractors are functions that transform imperfect random sources into nearly uniform random bits. Their original motivation is to bridge the gap between the uniform random bits required in standard applications (such as in randomized algorithms, distributed computing, and cryptography), and practical random sources which are almost always biased (either because of natural noise or adversarial information leakage). However the study of these objects has led to applications far beyond this motivation, in several different fields of computer science and combinatorics (e.g., coding theory, graph theory, and complexity theory).

As mentioned above, the inputs to a randomness extractor are usually imperfect randomness, which are modeled by the notion of general weak random sources with a certain amount of entropy.

Definition 1.1. The *min-entropy* of a random variable X is

$$H_\infty(X) = \min_{x \in \text{supp}(X)} \log_2(1/\Pr[X = x]).$$

For $X \in \{0, 1\}^n$, we call X an $(n, H_\infty(X))$ -source, and we say X has *entropy rate* $H_\infty(X)/n$.

An extensively studied model of randomness extractors is the so called *seeded extractors*, introduced by Nisan and Zuckerman [NZ96]. The inputs to a seeded extractor are a general weak random source and a short independent uniform random seed. The random seed is necessary here since it is well known that no deterministic extractor with one general weak source as input can exist. Seeded extractors have many applications in computer science, and we have the following formal definition.

Definition 1.2. (Seeded Extractor) A function $\text{Ext} : \{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^m$ is a (k, ϵ) -*extractor* if for every source X with min-entropy k and independent Y which is uniform on $\{0, 1\}^d$,

$$|\text{Ext}(X, Y) - U_m| \leq \epsilon.$$

If in addition we have $|(\text{Ext}(X, Y), Y) - (U_m, Y)| \leq \epsilon$ then we say it is a *strong* (k, ϵ) -*extractor*.

Through a long line of research, we now have explicit constructions of seeded extractors with almost optimal parameters (e.g., [LRVW03, GUV09, DW08, DKSS09]). In the last decade or so, the focus has shifted to several different but related models of randomness extractors, including seedless extractors and non-malleable extractors. The study of these topics has also been quite fruitful, leading to breakthroughs to several long standing open problems.

1.1 Seedless extractors

As the name suggests, a seedless extractor uses no uniform seed, and the only inputs are weak random sources. Here, again we have two different cases. In the first case, one puts additional restrictions on a single weak random source in order to allow possible extraction, thus obtaining deterministic extractors for special classes of (structured) sources. In the second case, the sources are still general weak random sources, but the extractor needs to use more than one sources. To make extraction possible, one typically assumes the input sources to the extractor are independent, and this kind of extractors are sometimes called independent source extractors.

Since the pioneering work of Chor and Goldreich [CG88], the study of independent source extractors has gained significant attention due to their close connections to explicit Ramsey graphs, and their applications in distributed computing and cryptography with general weak random sources [KLRZ08, KLR09]. The goal here is to give explicit constructions that match the probabilistic bound: an extractor for just two independent (n, k) sources with $k \geq \log n + O(1)$ that outputs $\Omega(k)$ bits with exponentially small (in k) error. Note that an explicit two-source extractor for such entropy (even with one bit output and constant error) will give an (strongly) explicit Ramsey graph on N vertices with no clique or independent set of size $O(\log N)$, solving an open problem proposed by Erdős [Erd47] in his seminal paper that inaugurated the probabilistic method.

While early progress on this problem has been quite slow, with the best known construction in almost 20 years only able to handle two independent (n, k) sources with $k > n/2$ [CG88], since 2004 there has been a long line of work [BIW04, BKS⁺05, Raz05, Bou05, Rao06, BRSW06, Li11, Li12b, Li13b, Li13a, Li15b, Coh15, CZ16, Li16, CS16, CL16, Coh16a, BADTS17, Coh17, Li17] introducing exciting new techniques to this problem. This line of work greatly improved the situation and led to a series of breakthroughs. Now we have three source extractors for entropy $k \geq \text{polylog}(n)$ that output $\Omega(k)$ bits with exponentially small error [Li15b], two-source extractors for entropy $k \geq \text{polylog}(n)$ that output $\Omega(k)$ bits with polynomially small error [CZ16, Li16, Mek15], and two-source extractors for entropy $k \geq O(\log n \log \log n)$ that output one bit with any constant error [Li17]. This also gives an explicit Ramsey graph on N vertices with no clique or independent set of size $(\log N)^{O(\log \log \log N)}$. Interestingly and somewhat surprisingly, the most recent progress which brought the entropy requirement close to optimal, has mainly benefited from the study of another kind of extractors, the so called *non-malleable extractors*, which we now describe below.

1.2 Non-malleable extractors

Non-malleable extractors are strengthening of standard extractors, where one requires that the output is close to uniform even given the output of the extractor on tampered inputs.

Definition 1.3 (Tampering Function). For any function $f : S \rightarrow S$, f has a fixed point at $s \in S$ if $f(s) = s$. We say f has no fixed points in $T \subseteq S$, if $f(t) \neq t$ for all $t \in T$. We say f has no fixed points if $f(s) \neq s$ for all $s \in S$. For any $n > 0$, let \mathcal{F}_n denote the set of all functions $f : \{0, 1\}^n \rightarrow \{0, 1\}^n$. Any subset of \mathcal{F}_n is a family of tampering functions.

Again, there are different models of non-malleable extractors. If the tampering acts on the seed of a seeded extractor, such extractors are called *seeded non-malleable extractors*, originally introduced by Dodis and Wichs in [DW09].

Definition 1.4 (Non-malleable extractor). A function $\text{snmExt} : \{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^m$ is a seeded non-malleable extractor for min-entropy k and error ϵ if the following holds: If X is a source on $\{0, 1\}^n$ with min-entropy k and $\mathcal{A} : \{0, 1\}^d \rightarrow \{0, 1\}^d$ is an arbitrary tampering function with no fixed points, then

$$|\text{snmExt}(X, U_d) \circ \text{snmExt}(X, \mathcal{A}(U_d)) \circ U_d - U_m \circ \text{snmExt}(X, \mathcal{A}(U_d)) \circ U_d| < \epsilon$$

where U_m is independent of U_d and X .

If the tampering acts on the sources of an independent source extractor, such extractors are called *seedless non-malleable extractors*, originally introduced by Cheraghchi and Guruswami [CG14b].

Definition 1.5 (Seedless Non-Malleable C -Source Extractor). A function $\text{nmExt} : (\{0, 1\}^n)^C \rightarrow \{0, 1\}^m$ is a (k, ϵ) -seedless non-malleable extractor for C independent sources, if it satisfies the following property: Let X_1, \dots, X_C be C independent (n, k) sources, and $f_1, \dots, f_C : \{0, 1\}^n \rightarrow \{0, 1\}^n$ be C arbitrary tampering functions such that there exists an f_i with no fixed points, then

$$|\text{nmExt}(X_1, \dots, X_C) \circ \text{nmExt}(f_1(X_1), \dots, f_C(X_2)) - U_m \circ \text{nmExt}(f_1(X_1), \dots, f_C(X_2))| < \epsilon.$$

Remark 1.6. The original definition of seedless non-malleable independent source extractors in [CG14b] is more general, in the sense that the tampering functions may have some fixed points. However, Cheraghchi and Guruswami [CG14b] showed that this is essentially equivalent to the definition above, up to a small loss in parameters. We present the details in Section 7.

Seeded non-malleable extractors and privacy amplification. As stated above, seeded non-malleable extractors were first introduced by Dodis and Wichs [DW09], to study a cryptographic problem known as *privacy amplification* [BBR88]. The problem considers the situation where two parties with local (non-shared) uniform random bits try to convert a shared secret weak random source \mathbf{X} into shared secret nearly uniform random bits. They do this by communicating through a channel, which is watched by an adversary with unlimited computational power. Standard strong seeded extractors provide very efficient protocols for a passive adversary (i.e., can only see the messages but cannot change them), but fail for an active adversary (i.e., can arbitrarily change, delete and reorder messages). In the latter case (which is the focus of this paper), the main goal is to design a protocol that uses as few number of interactions as possible, and achieves a shared uniform random string \mathbf{R} which has *entropy loss* (the difference between the length of the output and $H_\infty(\mathbf{X})$) as small as possible. Such a protocol is defined with a security parameter s , which means the probability that an active adversary can successfully make the two parties output two different strings without being detected is at most 2^{-s} . On the other hand, if the adversary remains passive, then the two parties should achieve shared secret random bits that are 2^{-s} -close to uniform. We refer the reader to [DLWZ14] for a formal definition.

A long line of work has been devoted to this problem [MW97, DKRS06, DW09, RW03, KR09, CKOR10, DLWZ14, CRS14, Li12a, Li12b, Li15a, CGL16, Coh16b, Coh16c, CL16, Coh16a, Coh17, Li17]. It is known that one round protocol can only exist when the entropy rate of \mathbf{X} is bigger than $1/2$, and the protocol has to incur a large entropy loss. When the entropy rate of \mathbf{X} is smaller than $1/2$, [DW09] showed that any protocol has to take at least two rounds with entropy loss at least $\Omega(s)$. Achieving a two-round protocol with entropy loss $O(s)$ for all possible security parameters s is thus the holy grail of this problem (note that s can be at most $\Omega(k)$ where $k = H_\infty(\mathbf{X})$).

While early works on this problem used various techniques, in [DW09], Dodis and Wichs introduced a major tool, the seeded non-malleable extractor defined above. They showed that two-round privacy amplification protocols with optimal entropy loss can be constructed using explicit seeded non-malleable extractors. Furthermore, non-malleable extractors exist when $k > 2m + 2 \log(1/\epsilon) + \log d + 6$ and $d > \log(n - k + 1) + 2 \log(1/\epsilon) + 5$. Since then, the study of non-malleable extractors has seen significant progress starting from the first explicit construction in [DLWZ14], with further connections to independent source extractors established in [Li12b, Li13b, CZ16]. Previous to this work, the best known seeded non-malleable extractor is due to the author [Li17], which works for entropy $k \geq O(\log n + \log(1/\epsilon) \log \log(1/\epsilon))$ and has seed length $d = O(\log n + \log(1/\epsilon) \log \log(1/\epsilon))$. Although quite close to optimal, the extra $O(\log \log(1/\epsilon))$ factor in the entropy requirement implies that by using this extractor, one can only get two-round privacy amplification protocols with

optimal entropy loss for security parameter up to $s = \Omega(k/\log k)$. This still falls short of achieving the holy grail, and may be problematic for some applications. For example, even if the shared weak source has slightly super-logarithmic entropy, the error of the protocol can still be sub-polynomially large; while ideally one can hope to get negligible error, which is important for other cryptographic applications based on this. The only previous protocol that can achieve security parameter up to $s = \Omega(k)$ is the work of [CKOR10], which has entropy loss $O(\log n + s)$ but also uses $O(\log n + s)$ rounds of interactions, much larger than 2. This also results in a total communication complexity of $O((\log n + s)^2)$ and requires the two parties' local random bits to be at least this long.

Seedless non-malleable extractors and non-malleable codes. Seedless non-malleable extractors were first introduced by Cheraghchi and Guruswami [CG14b] to study non-malleable codes [DPW10], a generalization of standard error correcting codes to handle a much larger class of attacks. Informally, a non-malleable code is defined w.r.t. a specific family of tampering functions \mathcal{F} . The code consists of a randomized encoding function E and a deterministic decoding function D , such that for any $f \in \mathcal{F}$, if a codeword $E(x)$ is modified into $f(E(x))$, then the decoded message $x' = D(f(E(x)))$ is either the original message x or a completely unrelated message. The formal definition is given in Section 7. In [DPW10], Dziembowski et. al showed that such codes can be used generally in tamper-resilient cryptography to protect the memory of a device.

Even with such generalization, non-malleable codes still cannot exist if \mathcal{F} is completely unrestricted. However, they do exist for many broad families of tampering functions. One of the most studied families of tampering functions is the so called *t-split-state* model. Here, a k -bit message x is encoded into a codeword with t parts y_1, \dots, y_t , each of length n . An adversary can then arbitrarily tamper with each y_i independently. In this case, the rate of the code is defined as $k/(tn)$.

This model arises naturally in many applications, typically when different parts of memory are used to store different parts of y_1, \dots, y_t . Such a code can also be viewed as a kind of “non-malleable secret sharing scheme”. The case of $t = 2$ is the most useful and interesting setting, since $t = 1$ corresponds to the case where \mathcal{F} is unrestricted. Again, there has been a lot of previous work on non-malleable codes in this model. In this paper we will focus on the information theoretic setting.

Dziembowski et. al [DPW10] first proved the existence of non-malleable codes in the split-state model. Cheraghchi and Guruswami [CG14a] showed that the optimal rate of such codes in the 2-split-state model is $1/2$. The first explicit construction appears in [DKO13], with later improvements appearing in [ADL14, Agg14, ADKO15], but all constructions only achieve rate $n^{-\Omega(1)}$.

Cheraghchi and Guruswami [CG14b] found a way to construct non-malleable codes in the t -split state model using sufficiently good non-malleable t -source extractors. Chattopadhyay and Zuckerman [CZ14] constructed the first seedless non-malleable, which works for 10 independent sources with entropy $(1-\gamma)n$. They further used this extractor to give a constant rate non-malleable code in the 10-split-state model. Subsequently, constructions of non-malleable two source extractors appeared in [CGL16] and [Li17], where both constructions work for min-entropy $k = (1-\gamma)n$ and output $\Omega(k)$ bits. The construction in [CGL16] has error $2^{-n^{\Omega(1)}}$ while the construction in [Li17] has error $2^{-\Omega(n/\log n)}$. Both can be used to give explicit non-malleable codes in the 2-split state model, where the former achieves rate $n^{-\Omega(1)}$ and the latter achieves $\Omega(\frac{1}{\log n})$. Very recently, a work by Kanukurthi et. al [KOS17] achieved constant rate in the 4-split state model, but the best construction in the 2-split state model still only achieves rate $\Omega(\frac{1}{\log n})$ [Li17].

As can be seen from the above discussions, extensive past research has established strong connections among the following 5 seemingly different problems: seeded non-malleable extractors, privacy

amplification protocols, independent source extractors (and explicit Ramsey graphs), non-malleable independent source extractors, and non-malleable codes (in the split state model). Furthermore, past research has brought each one of them close to optimal, but there remains a small gap to close and the motivation to close this gap remains strong. On the other hand, achieving this seems challenging and beyond the reach of the techniques developed so far.

1.3 Our Results

In this paper we introduce new techniques as an effort and the first step to close the gaps mentioned above. Our techniques lead to improvements to all the 5 problems discussed. We will first list our results here, and then give an informal overview of our techniques in the next section. Our first theorem gives explicit seeded non-malleable extractors which have optimal entropy requirement with respect to the error.

Theorem 1.7. *There exists a constant $C > 1$ such that for any constant $a \in \mathbb{N}, a \geq 2$, any $n, k \in \mathbb{N}$ and any $0 < \epsilon < 1$ with $k \geq C(\log \log n + a \log(1/\epsilon))$, there is an explicit construction of a strong seeded (k, ϵ) non-malleable extractor $\{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^m$ with $d = O(\log n) + \log(1/\epsilon)2^{O(a(\log \log(1/\epsilon))^{\frac{1}{a}})}$ and $m = \Omega(k)$.*

Note that this theorem provides a trade-off between the entropy requirement and the seed length. For example, if we take $a = 2$, then we need the source to have entropy $O(\log \log n + \log(1/\epsilon))$ while the seed length is $O(\log n) + 2^{O(\sqrt{\log \log(1/\epsilon)})} \log(1/\epsilon) = O(\log n) + \log^{1+o(1)}(1/\epsilon)$. By a recent result of Gur and Shinkar [GS18], the entropy requirement in our construction is asymptotically optimal. Combined with the protocol in [DW09], this gives the following theorem.

Theorem 1.8. *For any constant integer $a \geq 2$ there exists a constant $0 < \alpha < 1$ such that for any $n, k \in \mathbb{N}$ and security parameter $s \leq \alpha k$, there is an explicit two-round privacy amplification protocol with entropy loss $O(\log \log n + s)$, in the presence of an active adversary. The communication complexity of the protocol is $O(\log n) + s2^{O(a(\log s)^{\frac{1}{a}})}$.*

Note that our two-round protocol has optimal entropy loss for security parameter up to $s = \Omega(k)$, thus achieving the holy grail of this problem. Compared to the $O(\log n + s)$ -round protocol in [CKOR10], our protocol also has better dependence on n and significantly better communication complexity. We remark that the $O(\log \log n)$ term in both theorems is also the best possible (up to constant) if one wants to apply the two-round protocol in [DW09]. This is because the output of the non-malleable extractor is used in the second round as the key for a message authentication code (MAC) that authenticates the seed of a strong seeded extractor with security parameter s . Since the seed of the extractor uses at least $\Omega(\log n)$ bits, the MAC requires a key of length at least $\log \log n + s$. See [DW09] for more details.

We can also achieve smaller seed length while requiring slightly larger entropy.

Theorem 1.9. *There exists a constant $C > 1$ such that for any $n, k \in \mathbb{N}$ and $0 < \epsilon < 1$ with $k \geq C(\log \log n + \log(1/\epsilon) \log \log \log(1/\epsilon))$, there is an explicit construction of a strong seeded (k, ϵ) non-malleable extractor $\{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^m$ with $d = O(\log n + \log(1/\epsilon)(\log \log(1/\epsilon))^2)$ and $m = \Omega(k)$.*

Theorem 1.10. *There exists a constant $0 < \alpha < 1$ such that for any $n, k \in \mathbb{N}$ and security parameter $s \leq \alpha k / \log \log k$, there is an explicit two-round privacy amplification protocol with entropy*

loss $O(\log \log n + s)$, in the presence of an active adversary. The communication complexity of the protocol is $O(\log n + s \log^2 s)$.

Remark 1.11. In both Theorem 1.7 and Theorem 1.9, the dependence on the error ϵ in the seed length and entropy requirement can be switched. For example, in Theorem 1.7, we can also achieve $k \geq C \log \log n + \log(1/\epsilon) 2^{C \cdot a(\log \log(1/\epsilon))^{\frac{1}{a}}}$ and $d = O(\log n + a \log(1/\epsilon))$. In other words, our construction can be asymptotically optimal in either the seed length or the entropy requirement, but not in both.

We also have the following non-malleable two-source extractor and seeded non-malleable extractor.

Theorem 1.12. *There exists a constant $0 < \gamma < 1$ and a non-malleable two-source extractor for $(n, (1 - \gamma)n)$ sources with error $2^{-\Omega(n \log \log n / \log n)}$ and output length $\Omega(n)$.*

Theorem 1.13. *There is a constant $C > 0$ such that for any $\epsilon > 0$ and $n, k \in \mathbb{N}$ with $k \geq C(\log \log n + \frac{\log(1/\epsilon) \log \log(1/\epsilon)}{\log \log \log(1/\epsilon)})$, there is an explicit strong seeded non-malleable extractor for (n, k) sources with seed length $d = O(\log n + \frac{\log(1/\epsilon) \log \log(1/\epsilon)}{\log \log \log(1/\epsilon)})$, error ϵ and output length $\Omega(k)$.*

Combined with the techniques in [BADTS17], we obtain the following theorem which gives improved constructions of two-source extractors.

Theorem 1.14. *For every constant $\epsilon > 0$, there exists a constant $C > 1$ and an explicit two source extractor $\text{Ext} : (\{0, 1\}^n)^2 \rightarrow \{0, 1\}$ for entropy $k \geq C \frac{\log n \log \log n}{\log \log \log n}$ with error ϵ .*

As a corollary, we obtain the following improved constructions of Ramsey graphs.

Corollary 1.15. *For every large enough integer N there exists a (strongly) explicit construction of a K -Ramsey graph on N vertices with $K = (\log N)^{O(\frac{\log \log \log N}{\log \log \log \log N})}$.*

We can also efficiently sample uniformly from the pre-image of any given output of the extractor in Theorem 1.12. Combined with the connection in [CG14b], we obtain the following theorem.

Theorem 1.16. *For any $n \in \mathbb{N}$ there exists a non-malleable code with efficient encoder/decoder in the 2-split-state model with block length $2n$, rate $\Omega(\log \log n / \log n)$ and error $\epsilon = 2^{-\Omega(n \log \log n / \log n)}$.*

1.4 Overview of The Constructions and Techniques

We demonstrate our techniques here by an informal overview of our construction on non-malleable extractors. Throughout this section we will be mainly interested in the dependence of various parameters (e.g., seed length, entropy requirement) on the error ϵ , since this makes the presentation cleaner. The dependence on n comes from the alternating extraction between the seed and the source, thus the seed needs to have an $O(\log n)$ term while the source only needs an $O(\log \log n)$ term. We use letters with prime to denote the tampered version of random variables.

All recent constructions on non-malleable extractors essentially follow the same high level sketch: first obtain a small advice on $L = O(\log(1/\epsilon))$ bits such that with probability $1 - \epsilon$, the advice is different from its tampered version. Then, use the rest of the inputs, together with a correlation breaker with advice (informally introduced in [CGL16] and formally defined in [Coh16b]) to obtain the final output. There are several constructions of the correlation breaker, with the most efficient one using a non-malleable independence preserving merger (NIPM for short, introduced in [CS16]

and generalized in [CL16]). The NIPM takes an $L \times m$ random matrix V with $m = O(\log(1/\epsilon))$ and use the other inputs to merge it into one output. It has the property that if the matrix has one row which is uniform given the corresponding row in its tampered version² (which can be obtained from the advice and inputs), then the output is guaranteed to be uniform given the tampered output. From now on, we assume the inputs to the extractor are two independent sources X and Y (in the case of seeded non-malleable extractor, Y can be viewed as the seed).

Previously, the best construction of an NIPM is due to the author [Li17], which works roughly as follows. Suppose the matrix V is a deterministic function of the source X , then we first generate $\ell = \log L$ random variables (Y_1, \dots, Y_ℓ) from Y , such that each Y_i is close to uniform given the previous random variables and their tampered versions (i.e., $(Y_1, Y'_1, \dots, Y_{i-1}, Y'_{i-1})$). We call this property the *look-ahead* property. Next, we run a simpler merger for ℓ iterations, with each iteration using a new Y_i to merge every two consecutive rows in V , thus decreasing the number of rows by a factor of 2. We output the final matrix V which has one row. The property of the merger guarantees that in the end the output is uniform given its tampered version.

Let's turn to the entropy requirement. In this construction each Y_i needs to have at least $\Omega(\log(1/\epsilon))$ bits in order to ensure the error is at most ϵ , thus it is clear that Y needs to have entropy at least $O(\ell \log(1/\epsilon)) = O(\log(1/\epsilon) \log \log(1/\epsilon))$. However, it turns out that X also needs to have such entropy, for the following two reasons. First, in each iteration after we apply the simple merger, the length of each row in the matrix decreases by a constant factor (due to the entropy loss of any seeded extractor). Thus we cannot afford to just repeat the process for ℓ times since that would require the original row in V to have length at least $\text{polylog}(1/\epsilon)$, which implies the same entropy requirement for X . Instead, we again create ℓ random variables (X_1, \dots, X_ℓ) from X with the look-ahead property, and in each iteration after merging we use each row of the matrix to extract from a new X_i (using a standard seeded extractor, and possibly after first extracting from another new Y_i), in order to restore the length of the rows in the matrix. We need the look-ahead property in (X_1, \dots, X_ℓ) and (Y_1, \dots, Y_ℓ) so that after each iteration we can fix the previously used random variables and maintain the independence of X and Y , as well as the fact that the matrix is a deterministic function of X . Each X_i again needs at least $\Omega(\log(1/\epsilon))$ bits so this puts a lower bound on the entropy of X .

Second, in order to prepare the random variables (Y_1, \dots, Y_ℓ) , we in fact run an alternating extraction protocol between (part of) X and Y . This protocol lasts 2ℓ rounds between X and Y , and in each round either X or Y needs to spend $\Omega(\log(1/\epsilon))$ random bits. This again puts a lower bound of $O(\ell \log(1/\epsilon))$ on the entropy of X .

We remark that the above description is slightly different from the standard definition of an NIPM, where the only input besides the matrix V is Y . Indeed, in [Li17] it was presented as a correlation breaker. However, these two objects are actually similar, and for this paper it is more convenient to consider NIPMs with an additional input X , which is independent of Y but may be correlated with V . We will use this notion here and formally define it in Section 4.

Improved merger construction. We develop new techniques to break the above barriers, so that we can achieve essentially optimal entropy requirement in one of the sources (e.g., X). Our new techniques recycle the entropy in X , similar in spirit to what has been done in previous constructions of pseudorandom generators for small space computation [Nis92, NZ96].

²Sometimes we also require the other rows to be uniform, in order to make the construction simpler. This is the case of this paper, but we ignore the issue here for simplicity and clarity.

For the first problem, our key observation is that the random variables (X_1, \dots, X_ℓ) can be replaced by the original source X , as long as we have slightly more (e.g., 2ℓ) Y_i 's and they satisfy the look ahead property. To achieve this we crucially use the property that the NIPM only needs one row of V to be uniform given the corresponding row in its tampered version, and does not care about the dependence among the rows of V (in fact, they can have arbitrary dependence). Consider a particular iteration i in which we have just finished applying the simple merger. We can first fix all random variables $\{Y_j\}$ that have been used so far, and conditioned on this fixing we know that X and Y are still independent, and the matrix V is a deterministic function of X , which is independent of all random variables obtained from Y . To restore the length of each row in V , we use each row of V to first extract $O(\log(1/\epsilon))$ bits from Y_{j+1} , and then extract back from the original source X . Note that we only need to consider each row separately (since we don't care about the dependence among them). Assume row h in V has the property that V_h is uniform given V'_h . Since each random variable only has $O(\log(1/\epsilon))$ bits, as long as the entropy of X is $c \log(1/\epsilon)$ for a large enough constant $c > 1$, we can argue that conditioned on the fixing of (V_h, V'_h) , X still has entropy at least some $O(\log(1/\epsilon))$. On the other hand since V_h is uniform given V'_h , their corresponding outputs after extracting from (Y_{j+1}, Y'_{j+1}) will also preserve this independence; and conditioned on the fixing of (V_h, V'_h) , these outputs are deterministic functions of (Y, Y') , which are independent of (X, X') . Thus they can be used to extract back from (X, X') and preserve the independence. By standard properties of a strong seeded extractor, this holds even conditioned on the fixing of (Y_{j+1}, Y'_{j+1}) . Note that conditioned on the further fixing of (Y_{j+1}, Y'_{j+1}) , the new matrix is again a deterministic function of X , thus we can go into the next iteration. Therefore, by recycling the entropy in X , altogether we only need X to have entropy some $O(\log(1/\epsilon))$. In each iteration we use two new Y_i 's so we need roughly 2ℓ such random variables.

We note the following important difference between this approach and the previous approach: the previous approach actually also roughly preserves the independence between different rows in the matrix, since each time we use a new X_i to restore the length of each row, and X_i is uniform given all previously used random variables. In contrast, the current approach only preserves the independence between the corresponding row in V and V' , which is in fact all we need. Thus we can always use X to restore the length of each row, and this recycles the entropy in X .

However, we still need to address the second problem, where we need to generate the random variables $(Y_1, \dots, Y_{2\ell})$. One way to generate them with the look-ahead property, as we mentioned above, is to use an alternating extraction protocol, but this will require entropy roughly $O(\ell \log(1/\epsilon))$ from X . To solve this problem, we observe that there is another way to generate these random variables, which requires much less entropy from X . For simplicity assume that Y is uniform, we first take 2ℓ slices Y^i from Y , where Y^i has size $(2^i - 1)d$ for some $d = O(\log(1/\epsilon))$. This ensures that even conditioned on the fixing of $(Y^1, Y'^1, \dots, Y^{i-1}, Y'^{i-1})$, the (average) conditional min-entropy of Y_i is at least $(2^i - 1)d - 2 \cdot (2^{i-1} - 1)d = d$. Then, we can take $O(\log(1/\epsilon))$ uniform bits obtained from X , and use the *same* bits to extract Y_i from Y^i for every i . As long as we use a strong seeded extractor here, we are guaranteed that $(Y_1, \dots, Y_{2\ell})$ satisfy the look-ahead property; and moreover conditioned on the fixing of the $O(\log(1/\epsilon))$ bits from X , we have that $(Y_1, \dots, Y_{2\ell})$ is a deterministic function of Y . Note here again we only require entropy $O(\log(1/\epsilon))$ from X , and together with the approach described above this gives us a non-malleable extractor where X can have entropy $O(\log(1/\epsilon))$. However Y will need to have entropy at least $2^{2\ell} O(\log(1/\epsilon)) = O(\log^3(1/\epsilon))$.

To improve the entropy requirement of Y , we note that in the above approach, we only used part of X once to help obtaining the $\{Y^i\}$. Thus we have to use larger and larger slices of Y which

actually waste some entropy. Instead, we can use several parts of X , each with $O(\log(1/\epsilon))$ uniform bits. For example, say that we have obtained X^1 and X^2 , where each is uniform on some $O(\log(1/\epsilon))$ bits and X^2 is uniform even conditioned on the fixing of (X^1, X'^1) . We can now take some t slices $\{Y^i\}$ of Y , each of length $(2^i - 1) \cdot 2d$ for some parameters t, d . We first use X^1 to extract from each Y^i and obtain d uniform bits. Note that conditioned on the fixing of (X^1, X'^1) , these t random variables already satisfy the look-ahead property. Now for each of these d bits obtained from Y^i , we can apply the same process, i.e., we take some t slices of these d bits, each of length $(2^i - 1) \cdot O(\log(1/\epsilon))$ and then use X^2 to extract from each of them. This way we obtain t^2 random variables $\{Y_i\}$ that satisfy the look-ahead property. We can thus choose $t^2 = 2\ell$ which means $t = O(\sqrt{\ell})$. The entropy requirement of Y is roughly $(2^t - 1) \cdot (2^t - 1)O(\log(1/\epsilon)) = O(2^{2t} \log(1/\epsilon)) = 2^{O(\sqrt{\ell})} \log(1/\epsilon)$, while the entropy requirement for X is $2O(\log(1/\epsilon)) + O(\log(1/\epsilon)) = O(\log(1/\epsilon))$. This significantly improves the entropy requirement of Y .

We can repeat the previous process and use some a parts (X^1, \dots, X^a) obtained from X . As long as a is a constant integer, the entropy requirement for X will be $O(a \log(1/\epsilon)) = O(\log(1/\epsilon))$, while the entropy requirement of Y will be reduced to $2^{O(a\ell^{\frac{1}{a}})} \log(1/\epsilon) = 2^{O(a \log \log(1/\epsilon)^{\frac{1}{a}})} \log(1/\epsilon)$. To prepare the a parts of X , we perform an initial alternating extraction protocol between X and Y , which only needs entropy $O(a \log(1/\epsilon))$ from either of them. This gives Theorem 1.7. In the extreme case, we can also try to minimize the entropy requirement of Y . For this we can first create $\log \ell + 1 = \log \log \log(1/\epsilon) + O(1)$ X^i 's, and in each step use a new X^i to double the number of Y_i 's. This can be done by using the same X^i to do an alternating extraction of two rounds with each Y_i in parallel. Thus after $\log \ell + 1$ steps we obtain $(Y_1, \dots, Y_{2\ell})$. In this case X needs to have entropy $O(\log(1/\epsilon) \log \log \log(1/\epsilon))$. Ideally, we would want to claim that Y needs entropy $O(\log(1/\epsilon) \log \log(1/\epsilon))$, but due to technical reasons (each time the output shrinks by a constant factor) we can only show that this works as long as Y has entropy $O(\log(1/\epsilon)(\log \log(1/\epsilon))^2)$.³

The balanced case. Notice that in the above discussion, the entropy requirement for X and Y is unbalanced, in the sense that we can reduce one of them to be quite small, while the other is relatively large (in fact, larger than $O(\log(1/\epsilon) \log \log(1/\epsilon))$ as in previous construction [Li17]). For applications to two-source extractors and non-malleable codes, we need a balanced entropy requirement. Upon first look it does not seem that the new techniques we have introduced so far can achieve any improvement in this case, since in the above discussion we are still merging two rows of the matrix V in each step, and for this merging we need at least $O(\log(1/\epsilon))$ fresh random bits. Note that we need $\ell = \log L = \log \log(1/\epsilon)$ steps to finish the merging, thus it seems the total entropy requirement is at least $O(\log(1/\epsilon) \log \log(1/\epsilon))$.

Our key observation here is that we can apply the same idea of recycling entropy discussed above. Specifically, let us choose a parameter $t \in \mathbb{N}$ and we merge every t rows in the matrix V at each step, using some merger that we have developed above. For example, we can choose the merger which for merging t rows, requires X to have entropy $O(\log(1/\epsilon))$ and Y to have entropy $2^{O(\sqrt{\log t})} \log(1/\epsilon)$. This will take us $\frac{\log L}{\log t}$ steps to finish merging, and we will do it in the following way. First, we create $s = O(\frac{\log L}{\log t})$ random variables X_1, \dots, X_s that satisfy the look-ahead property. Then, in each step of the merging, we will use a new X_j . The X_j 's can be prepared by taking a small slice of both X and Y and do an alternating extraction protocol with $O(s)$ rounds, which consumes entropy $O(s \log(1/\epsilon)) = O(\frac{\log L}{\log t} \log(1/\epsilon))$ from both X and Y . However, in each step of

³The exponent 2 can be reduced to be arbitrarily close to $\log 3$.

the merging, we will *not* use fresh entropy from Y , but will recycle the entropy in Y . Note that by doing this, we are recycling the entropy in both X and Y . The recycling in X is done within each step of applying the small merger, while the recycling in Y is done between these steps.

To achieve this, consider a particular step i in the merging. Since we are using a new X_j in each step, we can fix all previous X_j 's that have been used and their tampered versions. Conditioned on this fixing, the matrix V obtained so far (and the tampered version V') is a deterministic function of Y , therefore independent of X . We now want to claim that conditioned on the random variable (V, V') , Y still has high entropy. If this is true then we can take a new X_{j+1} and apply a strong seeded extractor to Y using X_{j+1} as the seed, and the extracted random bits (which are deterministic functions of Y conditioned on the fixing of X_{j+1}) can be used for merging in the next step. Also note that to apply the merger, we can take yet another new X_{j+2} and use each row of V to extract from X_{j+2} and create a matrix W . Conditioned on the fixing of (V, V') , we have that (W, W') is a deterministic function of (X, X') and therefore independent of (Y, Y') . Moreover the independence between corresponding rows in (V, V') is preserved in (W, W') (i.e., there is also a row in W that is uniform given the corresponding row in W'). Thus now we can indeed apply the merger again to W and the extracted random bits from Y , possibly together with a new X_{j+3} . Again, this is similar in spirit to what has been done in previous constructions of pseudorandom generators for small space computation [Nis92, NZ96].

The above idea indeed works, except for the following subtle point: the computation of the merger is actually not a small space computation. Indeed, the matrix V can have many rows and the size of V can be potentially larger than the entropy of Y (in fact, it will definitely be larger than the entropy of Y , unless we are willing to afford entropy $O(\log^2(1/\epsilon))$ in Y), at least in the first several steps of merging when the number of rows in V is still large. To get around this, we again use the property that the only thing we need from the matrix V is that one row in V is independent of the corresponding row in V' (we call this the good row), and between different rows we can allow arbitrary dependence. Therefore, since we are merging t rows in each step, we only need to condition on the fixing of these t rows (and their tampered versions). This will ensure that if originally there is a good row in these t rows, then after merging the output is still a good row in the new matrix. The dependence between different rows in the new matrix could change, but that does not matter. Thus, we only need the entropy of Y to be roughly $O(t \log(1/\epsilon)) + 2^{O(\sqrt{\log t})} \log(1/\epsilon) + O(\frac{\log L}{\log t} \log(1/\epsilon)) = O(t \log(1/\epsilon) + \frac{\log L}{\log t} \log(1/\epsilon))$ since we will maintain the length of each row in V to be $O(\log(1/\epsilon))$.

Now the problem is just to pick an appropriate t to make the entropy requirements for X and Y balanced. A simple calculation shows that by choosing $t = \frac{\log L}{\log \log L}$, both X and Y only need entropy $O(\frac{\log L}{\log \log L} \log(1/\epsilon)) = O(\frac{\log(1/\epsilon) \log \log(1/\epsilon)}{\log \log \log(1/\epsilon)})$. By the connections in [Li17, BADTS17, CG14b], this dependence will translate into two source extractors for entropy $O(\frac{\log n \log \log n}{\log \log \log n})$, non-malleable two-source extractors for entropy $(1 - \gamma)n$ with some constant $\gamma > 0$ and error $2^{-\Omega(n \log \log n / \log n)}$, and non-malleable codes in the 2-split state model with rate $\Omega(\log \log n / \log n)$.

Organization. The rest of the paper is organized as follows. We give some preliminaries in Section 2, and define alternating extraction in Section 3. We present independence preserving mergers in Section 4, correlation breakers in Section 5, non-malleable extractors in Section 6, and non-malleable codes in Section 7. Finally we conclude with some open problems in Section 8.

2 Preliminaries

We often use capital letters for random variables and corresponding small letters for their instantiations. Let $|S|$ denote the cardinality of the set S . For ℓ a positive integer, U_ℓ denotes the uniform distribution on $\{0, 1\}^\ell$. When used as a component in a vector, each U_ℓ is assumed independent of the other components. When we have adversarial tampering, we use letters with prime to denote the tampered version of random variables. All logarithms are to the base 2.

2.1 Probability Distributions

Definition 2.1 (statistical distance). Let W and Z be two distributions on a set S . Their *statistical distance* (variation distance) is

$$\Delta(W, Z) \stackrel{\text{def}}{=} \max_{T \subseteq S} (|W(T) - Z(T)|) = \frac{1}{2} \sum_{s \in S} |W(s) - Z(s)|.$$

We say W is ε -close to Z , denoted $W \approx_\varepsilon Z$, if $\Delta(W, Z) \leq \varepsilon$. For a distribution D on a set S and a function $h : S \rightarrow T$, let $h(D)$ denote the distribution on T induced by choosing x according to D and outputting $h(x)$.

Lemma 2.2. For any function α and two random variables A, B , we have $\Delta(\alpha(A), \alpha(B)) \leq \Delta(A, B)$.

2.2 Average Conditional Min Entropy

Definition 2.3. The *average conditional min-entropy* is defined as

$$\begin{aligned} \tilde{H}_\infty(X|W) &= -\log \left(\mathbb{E}_{w \leftarrow W} \left[\max_x \Pr[X = x | W = w] \right] \right) \\ &= -\log \left(\mathbb{E}_{w \leftarrow W} \left[2^{-H_\infty(X|W=w)} \right] \right). \end{aligned}$$

Lemma 2.4 ([DORS08]). For any $s > 0$, $\Pr_{w \leftarrow W} [H_\infty(X|W = w) \geq \tilde{H}_\infty(X|W) - s] \geq 1 - 2^{-s}$.

Lemma 2.5 ([DORS08]). If a random variable B has at most 2^ℓ possible values, then $\tilde{H}_\infty(A|B) \geq H_\infty(A) - \ell$.

2.3 Prerequisites from Previous Work

Sometimes it is convenient to talk about average case seeded extractors, where the source X has average conditional min-entropy $\tilde{H}_\infty(X|Z) \geq k$ and the output of the extractor should be uniform given Z as well. The following lemma is proved in [DORS08].

Lemma 2.6. [DORS08] For any $\delta > 0$, if Ext is a (k, ϵ) extractor then it is also a $(k + \log(1/\delta), \epsilon + \delta)$ average case extractor.

For a strong seeded extractor with optimal parameters, we use the following extractor constructed in [GUV09].

Theorem 2.7 ([GUV09]). *For every constant $\alpha > 0$, there exists a constant $\beta > 0$ such that for all positive integers n, k and any $\epsilon > 2^{-\beta k}$, there is an explicit construction of a strong (k, ϵ) -extractor $\text{Ext} : \{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^m$ with $d = O(\log n + \log(1/\epsilon))$ and $m \geq (1 - \alpha)k$. The same statement also holds for a strong average case extractor.*

Theorem 2.8 ([CG88]). *For every $0 < m < n$ there is an explicit two-source extractor $\text{IP} : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}^m$ based on the inner product function, such that if X, Y are two independent (n, k_1) and (n, k_2) sources respectively, then*

$$(\text{IP}(X, Y), X) \approx_\epsilon (U_m, X) \text{ and } (\text{IP}(X, Y), Y) \approx_\epsilon (U_m, Y),$$

where $\epsilon = 2^{-\frac{k_1+k_2-n-m-1}{2}}$.

The following standard lemma about conditional min-entropy is implicit in [NZ96] and explicit in [MW97].

Lemma 2.9 ([MW97]). *Let X and Y be random variables and let \mathcal{Y} denote the range of Y . Then for all $\epsilon > 0$, one has*

$$\Pr_Y \left[H_\infty(X|Y = y) \geq H_\infty(X) - \log |\mathcal{Y}| - \log \left(\frac{1}{\epsilon} \right) \right] \geq 1 - \epsilon.$$

We also need the following lemma.

Lemma 2.10. [Li13a] *Let (X, Y) be a joint distribution such that X has range \mathcal{X} and Y has range \mathcal{Y} . Assume that there is another random variable X' with the same range as X such that $|X - X'| = \epsilon$. Then there exists a joint distribution (X', Y) such that $|(X, Y) - (X', Y)| = \epsilon$.*

3 Alternating Extraction

Our constructions use the following alternating extraction protocol as a key ingredient. Alternating extraction was first introduced in [DP07], and has now become an important tool in constructions related to extractors.

Definition 3.1. (Alternating Extraction) Assume that we have two parties, Quentin and Wendy. Quentin has a source Q , Wendy has a source W . Also assume that Quentin has a uniform random seed S_1 (which may be correlated with Q). Suppose that (Q, S_1) is kept secret from Wendy and W is kept secret from Quentin. Let $\text{Ext}_q, \text{Ext}_w$ be strong seeded extractors with optimal parameters, such as that in Theorem 2.7. Let r, s be two integer parameters for the protocol. For some integer parameter $\ell > 0$, the *alternating extraction protocol* is an interactive process between Quentin and Wendy that runs in ℓ steps.

In the first step, Quentin sends S_1 to Wendy, Wendy computes $R_1 = \text{Ext}_w(W, S_1)$. She sends R_1 to Quentin and Quentin computes $S_2 = \text{Ext}_q(Q, R_1)$. In this step R_1, S_2 each outputs r and s bits respectively. In each subsequent step i , Quentin sends S_i to Wendy, Wendy computes $R_i = \text{Ext}_w(W, S_i)$. She replies R_i to Quentin and Quentin computes $S_{i+1} = \text{Ext}_q(Q, R_i)$. In step i , R_i, S_{i+1} each outputs r and s bits respectively. Therefore, this process produces the following sequence:

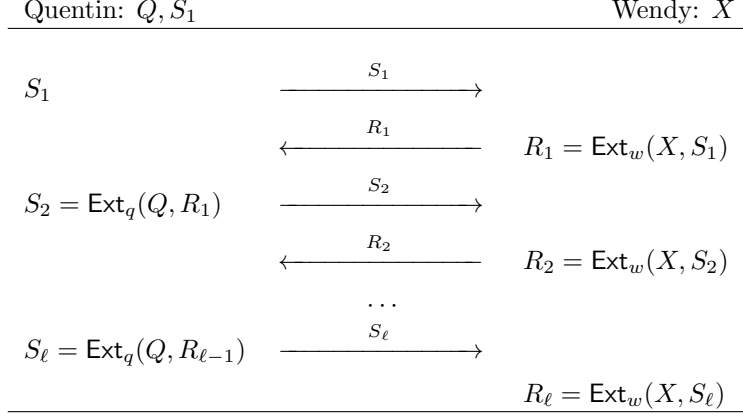


Figure 1: Alternating Extraction.

$$S_1, R_1 = \text{Ext}_w(W, S_1), S_2 = \text{Ext}_q(Q, R_1), \dots,$$

$$S_\ell = \text{Ext}_q(Q, R_{\ell-1}), R_\ell = \text{Ext}_w(W, S_\ell).$$

The output of an alternating extraction protocol is often described as a *look-ahead extractor*, defined as follows. Let $Y = (Q, S_1)$ be a seed, the look-ahead extractor is defined as

$$\text{laExt}(W, Y) = \text{laExt}(W, (Q, S_1)) \stackrel{\text{def}}{=} R_1, \dots, R_\ell.$$

The following lemma is a special case of Lemma 6.5 in [CGL16].

Lemma 3.2. *Let W be an (n_w, k_w) -source and W' be a random variable on $\{0, 1\}^{n_w}$ that is arbitrarily correlated with W . Let $Y = (Q, S_1)$ such that Q is a (n_q, k_q) -source, S_1 is a uniform string on s bits, and $Y' = (Q', S'_1)$ be a random variable arbitrarily correlated with Y , where Q' and S'_1 are random variables on n_q bits and s bits respectively. Let $\text{Ext}_q, \text{Ext}_w$ be strong seeded extractors that extract s and r bits from sources with min-entropy k with error ϵ and seed length $d \leq \min\{r, s\}$. Suppose (Y, Y') is independent of (W, W') , $k_q \geq k + 2(\ell - 1)s + 2 \log(\frac{1}{\epsilon})$, and $k_w \geq k + 2(\ell - 1)r + 2 \log(\frac{1}{\epsilon})$. Let laExt be the look-ahead extractor defined above using $\text{Ext}_q, \text{Ext}_w$, and $(R_1, \dots, R_\ell) = \text{laExt}(W, Y)$, $(R'_1, \dots, R'_\ell) = \text{laExt}(W', Y')$. Then for any $0 \leq j \leq \ell - 1$, we have*

$$(Y, Y', \{R_1, R'_1, \dots, R_j, R'_j\}, R_{j+1})$$

$$\approx_{\epsilon_1} (Y, Y', \{R_1, R'_1, \dots, R_j, R'_j\}, U_r),$$

where $\epsilon_1 = O(\ell\epsilon)$.

4 Non-Malleable Independence Preserving Merger

We now describe the notion of *non-malleable independence preserving merger*, introduced in [CL16] based on the notion of independence preserving merger introduced in [CS16].

Definition 4.1. A (L, d', ε) -NIPM $: \{0, 1\}^{Lm} \times \{0, 1\}^d \rightarrow \{0, 1\}^{m_1}$ satisfies the following property. Suppose

- \mathbf{X}, \mathbf{X}' are random variables, each supported on boolean $L \times m$ matrices s.t for any $i \in [L]$, $\mathbf{X}_i = U_m$,
- $\{\mathbf{Y}, \mathbf{Y}'\}$ is independent of $\{\mathbf{X}, \mathbf{X}'\}$, s.t \mathbf{Y}, \mathbf{Y}' are each supported on $\{0, 1\}^d$ and $H_\infty(\mathbf{Y}) \geq d'$,
- there exists an $h \in [L]$ such that $(\mathbf{X}_h, \mathbf{X}'_h) = (U_m, \mathbf{X}'_h)$,

then

$$|(L, d', \varepsilon)\text{-NIPM}(\mathbf{X}, \mathbf{Y}), (L, d', \varepsilon)\text{-NIPM}(\mathbf{X}', \mathbf{Y}') - U_{m_1}, (L, d', \varepsilon)\text{-NIPM}(\mathbf{X}', \mathbf{Y}')| \leq \varepsilon.$$

We have the following construction and theorem.

L -Alternating Extraction We extend the previous alternating extraction protocol by letting Quentin have access to L sources Q_1, \dots, Q_L (instead of just Q) which have the same length. Now in the i 'th round of the protocol, he uses Q_i to produce the r.v $S_i = \text{Ext}_q(Q_i, R_i)$. More formally, the following sequence of r.v's is generated: $S_1, R_1 = \text{Ext}_w(W, S_1), S_2 = \text{Ext}_q(Q_2, R_1), \dots, R_{L-1} = \text{Ext}_w(W, S_{L-1}), S_L = \text{Ext}_q(Q_L, R_{L-1})$.

The NIPM is now constructed as follows. Let S_1 be a slice of \mathbf{X}_1 with length $O(\log(d/\varepsilon))$, then run the L -alternating extraction described above with $(Q_1, \dots, Q_L) = (\mathbf{X}_1, \dots, \mathbf{X}_L)$ and $W = \mathbf{Y}$. Finally output S_L .

Theorem 4.2 ([CL16]). *There exists a constant $c > 0$ such that for all integers $m, d, d', L > 0$ and any $\varepsilon > 0$, with $m \geq 4cL \log(d/\varepsilon)$, $d' \geq 4cL \log(m/\varepsilon)$, the above construction NIPM $: (\{0, 1\}^m)^\ell \times \{0, 1\}^d \rightarrow \{0, 1\}^{m_1}$ has output length $m_1 \geq 0.2m$, such that if the following conditions hold:*

- \mathbf{X}, \mathbf{X}' are random variables, each supported on boolean $L \times m$ matrices s.t for any $i \in [L]$, $\mathbf{X}_i = U_m$,
- $\{\mathbf{Y}, \mathbf{Y}'\}$ is independent of $\{\mathbf{X}, \mathbf{X}'\}$, s.t \mathbf{Y}, \mathbf{Y}' are each supported on $\{0, 1\}^d$ and $H_\infty(\mathbf{Y}) \geq d'$,
- there exists an $h \in [L]$ such that $(\mathbf{X}_h, \mathbf{X}'_h) = (U_m, \mathbf{X}'_h)$,

then

$$|\text{NIPM}(\mathbf{X}, \mathbf{Y}), \text{NIPM}(\mathbf{X}', \mathbf{Y}'), \mathbf{Y}, \mathbf{Y}' - U_{m_1}, \text{NIPM}(\mathbf{X}', \mathbf{Y}'), \mathbf{Y}, \mathbf{Y}'| \leq L\varepsilon.$$

It is sometimes more convenient to consider NIPMs which use an additional source X in the computation. We generalize the above definition as follows.

Definition 4.3. A (L, d, d', ε) -NIPM $: \{0, 1\}^{Lm} \times \{0, 1\}^d \times \{0, 1\}^{d'} \rightarrow \{0, 1\}^{m_1}$ satisfies the following property. Suppose

- V, V' are random variables, each supported on boolean $L \times m$ matrices s.t for any $i \in [L]$, $V_i = U_m$,
- there exists an $h \in [L]$ such that $(V_h, V'_h) = (U_m, V'_h)$,

- \mathbf{X}, \mathbf{X}' are random variables, each supported on d bits, such that \mathbf{X} is uniform conditioned on (V, V') ,
- $(\mathbf{Y}, \mathbf{Y}')$ is independent of $(V, V', \mathbf{X}, \mathbf{X}')$, s.t \mathbf{Y}, \mathbf{Y}' are each supported on $\{0, 1\}^{d'}$ and \mathbf{Y} is uniform,

If the function is an NIPM that is strong in \mathbf{Y} then

$$|(L, d, d', \varepsilon)\text{-NIPM}(V, \mathbf{X}, \mathbf{Y}), (L, d, d', \varepsilon)\text{-NIPM}(V', \mathbf{X}', \mathbf{Y}'), \mathbf{Y}, \mathbf{Y}' - U_{m_1}, (L, d, d', \varepsilon)\text{-NIPM}(V', \mathbf{X}', \mathbf{Y}'), \mathbf{Y}, \mathbf{Y}'| \leq \varepsilon.$$

If the function is an NIPM that is strong in \mathbf{X} then

$$|(L, d, d', \varepsilon)\text{-NIPM}(V, \mathbf{X}, \mathbf{Y}), (L, d, d', \varepsilon)\text{-NIPM}(V', \mathbf{X}', \mathbf{Y}'), \mathbf{X}, \mathbf{X}' - U_{m_1}, (L, d, d', \varepsilon)\text{-NIPM}(V', \mathbf{X}', \mathbf{Y}'), \mathbf{X}, \mathbf{X}'| \leq \varepsilon.$$

We will now use the above construction to give another NIPM, which recycles the entropy. Specifically, we have the following construction.

Construction 4.4. Asymmetric NIPM.

Inputs:

- $L, m, n, d \in \mathbb{N}$ and an error parameter $\varepsilon > 0$ such that $m \geq c \log(d/\varepsilon)$ and $d \geq c \log(n/\varepsilon)$ for some constant $c > 1$.
- A random variable V supported on a boolean $L \times m$ matrix.
- An $(n, 6m)$ source \mathbf{X} .
- Random variables $\mathbf{Y}_1, \dots, \mathbf{Y}_\ell$ where $\ell = \log L$ and each \mathbf{Y}_i is supported on $\{0, 1\}^d$.

Output: a random variable $\mathbf{W} \in \{0, 1\}^m$.

Let $V^0 = V$. For $i = 1$ to $\log L$ do the following.

1. Take a slice \mathbf{Y}_i^1 of \mathbf{Y}_i with length $d/3$. Merge every two rows of V^{i-1} , using \mathbf{Y}_i^1 and the NIPM from Theorem 4.2. That is, for every $j \leq t/2$ where t is the current number of rows in V^{i-1} (initially $t = L$), compute $\overline{V_j^{i-1}} = \text{NIPM}((V_{2j-1}^{i-1}, V_{2j}^{i-1}), \mathbf{Y}_i^1)$.
2. For every $j \leq t/2$, compute $\overline{\mathbf{Y}_{ij}} = \text{Ext}_1(\mathbf{Y}_i, \overline{V_j^{i-1}})$, where Ext_1 is the extractor in Theorem 2.7 and output $d/4$ bits.
3. For every $i \leq t/2$, compute $\widetilde{V_j^{i-1}} = \text{Ext}_2(\mathbf{X}, \overline{\mathbf{Y}_{ij}})$, where Ext_2 is the extractor in Theorem 2.7 and output m bits.
4. Let V^i with the concatenation of $\widetilde{V_j^{i-1}}, j = 1, \dots, t/2$. Note that the number of rows in V^i has decreased by a factor of 2.

Finally output $\mathbf{W} = V^{\log L}$.

Lemma 4.5. *There is a constant $c > 1$ such that suppose we have the following random variables:*

- V, V' , each supported on a boolean $L \times m$ matrix s.t for any $i \in [L]$, $V_i = U_m$. In addition, there exists an $h \in [L]$ such that $(V_h, V'_h) = (U_m, V'_h)$.
- \mathbf{X}, \mathbf{X}' where \mathbf{X} is an $(n, 6m)$ source.
- Random variables $(\mathbf{Y}_1, \mathbf{Y}'_1), \dots, (\mathbf{Y}_\ell, \mathbf{Y}'_\ell)$ obtained from \mathbf{Y}, \mathbf{Y}' deterministically, where $\ell = \log L$. These random variables satisfy the following look-ahead condition: $\forall j < \ell$, we have

$$(\mathbf{Y}_j, \mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{j-1}, \mathbf{Y}'_{j-1}) = (U_d, \mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{j-1}, \mathbf{Y}'_{j-1}).$$

In addition, $(V, V', \mathbf{X}, \mathbf{X}')$ is independent of $(\mathbf{Y}, \mathbf{Y}')$.

Let \mathbf{W} be the output of the NIPM on $(V, \mathbf{X}, \mathbf{Y}_1, \dots, \mathbf{Y}_\ell)$ and \mathbf{W}' be the output of the NIPM on $(V', \mathbf{X}', \mathbf{Y}'_1, \dots, \mathbf{Y}'_\ell)$. Then

$$(\mathbf{W}, \mathbf{W}', \mathbf{Y}, \mathbf{Y}') \approx_{O(L\epsilon)} (U_m, \mathbf{W}', \mathbf{Y}, \mathbf{Y}').$$

Proof. We use induction to show the following claim.

Claim 4.6. *For every $0 \leq i \leq \ell = \log L$, the following holds after step i .*

- V^i, V'^i are each supported on boolean $(t = L/2^i) \times m$ matrices s.t for any $j \in [t]$, $(V_j^i, \mathbf{Y}, \mathbf{Y}') \approx_{\epsilon_j} (U_m, \mathbf{Y}, \mathbf{Y}')$. In addition, there exists an $h \in [t]$ such that $(V_h^i, V'^i, \mathbf{Y}, \mathbf{Y}') \approx_{\epsilon_i} (U_m, V_h^i, \mathbf{Y}, \mathbf{Y}')$. Here ϵ_i is the error after step i which satisfies that $\epsilon_0 = 0$ and $\epsilon_{i+1} \leq 2\epsilon_i + 4\epsilon$.
- Conditioned on the fixing of $\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_j, \mathbf{Y}'_j$, each of V^i and V'^i is a deterministic function of $V, V', \mathbf{X}, \mathbf{X}'$.

For the base case of $i = 0$, the claim clearly holds. Now assume that the claim holds for i , we show that it holds for $i + 1$.

We first fix $\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_i, \mathbf{Y}'_i$. By the induction hypothesis, conditioned on the fixing of these random variables, each of V^i and V'^i is a deterministic function of $V, V', \mathbf{X}, \mathbf{X}'$, and thus independent of $(\mathbf{Y}_{i+1}, \mathbf{Y}'_{i+1})$. We only consider the row $h \in [t]$ such that $(V_h, V'_h) \approx_{4 \cdot 2^i \epsilon} (U_m, V'_h)$, since the analysis for the rest of the rows are similar and simpler.

First we ignore the error ϵ_i . By Theorem 4.2, and note that we are merging every two rows at one step, we can choose a suitable constant $c > 1$ in the construction such that

$$(\overline{V_{h'}^i}, \overline{V_{h'}^i}, \mathbf{Y}_{i+1}^1, \mathbf{Y}_{i+1}^1) \approx_{2\epsilon} (U_{m_1}, \overline{V_{h'}^i}, \mathbf{Y}_{i+1}^1, \mathbf{Y}_{i+1}^1),$$

where $h' = \lceil \frac{h}{2} \rceil$ and $m_1 = 0.2m$. We now fix $(\mathbf{Y}_{i+1}^1, \mathbf{Y}_{i+1}^1)$. Note that conditioned on the fixing, \mathbf{Y}_{i+1} still has average conditional min-entropy at least $d - d/3 = 2d/3$ and is independent of $(\overline{V_{h'}^i}, \overline{V_{h'}^i})$. Now we can first fix $\overline{V_{h'}^i}$ and then $\overline{V_{h'}^i}$. Note that conditioned on this fixing, $\overline{V_{h'}^i}$ is still (close to) uniform and the average conditional min-entropy of \mathbf{Y}_{i+1} is at least $2d/3 - d/4 > d/3$. Thus as long as c is large enough, by Theorem 2.7 we have that

$$(\overline{\mathbf{Y}_{ih'}}, \overline{V_{h'}^i}) \approx_\epsilon (U_{d/4}, \overline{V_{h'}^i}).$$

We now further fix $\overline{V_{h'}^i}$. Note that conditioned on this fixing, $\overline{\mathbf{Y}_{ih'}}$ is still (close to) uniform. Moreover conditioned on all the random variables we have fixed, $\overline{\mathbf{Y}_{ih'}}$ is a deterministic function of $\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{i+1}, \mathbf{Y}'_{i+1}$ and thus independent of \mathbf{X}, \mathbf{X}' . Also conditioned on all the random variables we have fixed, the average conditional min-entropy of \mathbf{X} is at least $6m - 2m_1 > 5m$.

We can now further fix $\widetilde{V_{h'}^i}$, which is a deterministic function of \mathbf{X}' . Conditioned on this fixing the independence of random variables still holds, while the average conditional min-entropy of \mathbf{X} is at least $5m - m = 4m$. Therefore by Theorem 2.7 we have that

$$(\widetilde{V_{h'}^i}, \overline{\mathbf{Y}_{ih'}}) \approx_\epsilon (U_m, \overline{\mathbf{Y}_{ih'}}).$$

Since we have already fixed $\overline{\mathbf{Y}'_{ih'}}$ and $\widetilde{V_{h'}^i}$, and note that conditioned on this fixing, $(\mathbf{Y}, \mathbf{Y}')$ are independent of $\widetilde{V_{h'}^i}$ which is a deterministic function of \mathbf{X} , we also have that

$$(\widetilde{V_{h'}^i}, \overline{\mathbf{X}'_{ih'}}, \mathbf{Y}, \mathbf{Y}') \approx_\epsilon (U_m, \overline{\mathbf{X}'_{ih'}}, \mathbf{Y}, \mathbf{Y}').$$

Adding back all the errors we get that there exists an $h' \in [t]$ such that

$$(\overline{\mathbf{X}_{h'}}, \widetilde{V_{h'}^i}, \mathbf{Y}, \mathbf{Y}') \approx_{\epsilon_{i+1}} (U_m, \widetilde{V_{h'}^i}, \mathbf{Y}, \mathbf{Y}'),$$

where $\epsilon_{i+1} \leq 2\epsilon_i + 4\epsilon$. Furthermore, it is clear that conditioned on the fixing of $\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{i+1}, \mathbf{Y}'_{i+1}$, each of V^{i+1} and V'^{i+1} is a deterministic function of $V, V', \mathbf{X}, \mathbf{X}'$.

We can now estimate the final error to be $\epsilon_\ell \leq 4(\sum_{i=1}^\ell 2^i \epsilon) = O(L\epsilon)$. Finally, when the number or rows in V^i decreases to 1 after step ℓ , the output $\mathbf{W} = V^{\log L}$ satisfies the conclusion of the lemma. \square

We will now construct another NIPM. First we need the following lemma.

Lemma 4.7. *For any constant $a \in \mathbb{N}$, any $\ell, s \in \mathbb{N}$ and any $\epsilon > 0$ there exists an explicit function $\text{Conv}_a : \{0, 1\}^n \times \{0, 1\}^{a \cdot d} \rightarrow \{0, 1\}^{\ell \cdot s}$ with $d = O(\log(n/\epsilon))$ and $n = 2^{O(a \cdot \ell^{\frac{1}{2}})} \cdot s$ such that the following holds. Let (Y, Y') be two random variables each on n bits, and Y is uniform. Let $(X = (X_1, \dots, X_a), X' = (X'_1, \dots, X'_a))$ be random variables each on $a \cdot d$ bits, where each X_i and X'_i is on d bits. Further assume that (X, X') satisfies the following look-ahead property: $\forall i \in [a]$, we have*

$$(X_i, X_1, X'_1, \dots, X_{i-1}, X'_{i-1}) = (U_d, X_1, X'_1, \dots, X_{i-1}, X'_{i-1}).$$

Let $(W_1, \dots, W_\ell) = \text{Conv}_a(Y, X)$ and $(W'_1, \dots, W'_\ell) = \text{Conv}_a(Y', X')$. Then we have

$$(X, X', W_1, W'_1, \dots, W_\ell, W'_\ell) \approx_{O(\ell\epsilon)} (X, X', U_s, W'_1, \dots, U_s, W'_\ell),$$

where each U_s is independent of previous random variables but may depend on later random variables.

Proof. We will prove the lemma by induction on a . For the base case $a = 1$, consider the following construction. For $j = 1, \dots, \ell$, let Y_j be a slice of Y with length $(2^j - 1) \cdot 2s$ (this is possible since the total entropy required is at most $2^\ell \cdot 2s$), and compute $W_j = \text{Ext}(Y_j, X_1)$. Note that for any $j \in [\ell]$, conditioned on the fixing of $Y_1, Y'_1, \dots, Y_{j-1}, Y'_{j-1}$, the average conditional min-entropy of Y_j is at least $(2^j - 1) \cdot 2s - 2(2^{j-1} - 1) \cdot 2s = 2s$. Thus by Theorem 2.7 we have that

$$(W_j, Y_1, Y'_1, \dots, Y_{j-1}, Y'_{j-1}, X, X') \approx_\epsilon (U_s, Y_1, Y'_1, \dots, Y_{j-1}, Y'_{j-1}, X, X').$$

Since $(W_1, W'_1, \dots, W_{j-1}, W'_{j-1})$ is a deterministic function of $(Y_1, Y'_1, \dots, Y_{j-1}, Y'_{j-1})$ and (X, X') , we also have that

$$(W_j, W_1, W'_1, \dots, W_{j-1}, W'_{j-1}, X, X') \approx_\epsilon (U_s, W_1, W'_1, \dots, W_{j-1}, W'_{j-1}, W, W').$$

By adding all the errors the statement of the lemma holds.

Now assume that the lemma holds for a , we will construct another function Conv_{a+1} for the case of $a + 1$. First choose a parameter $t \in \mathbb{N}$ to be decided later. For $j = 1, \dots, \ell/t$, let Y_j be a slice of Y with length $(2^j - 1) \cdot 2m$, where m is the length of Y (i.e., n) for Conv_a when choosing $\ell = t$. Thus we have $m = 2^{O(a \cdot t^{\frac{1}{a}})} \cdot s$. Now, for every j we first use X_1 to compute $\hat{W}_j = \text{Ext}(Y_j, X_1)$ and output m bits, then compute $(\hat{W}_{1j}, \dots, \hat{W}_{tj}) = \text{Conv}_a(\hat{W}_j, X_2, \dots, X_{a+1})$. The final outputs are obtained by combining all the $\{\hat{W}_{ij}\}$ in sequence.

Note that by the same argument as above, we have that

$$(X_1, X'_1, \hat{W}_1, \hat{W}'_1, \dots, \hat{W}_{\ell/t}, \hat{W}'_{\ell/t}) \approx_{O(\frac{\ell}{t}\epsilon)} (X_1, X'_1, U_m, \hat{W}'_1, \dots, U_m, \hat{W}'_{\ell/t}).$$

Now we can fix (X_1, X'_1) . Note that conditioned on the fixing, $(\hat{W}_1, \hat{W}'_1, \dots, \hat{W}_{\ell/t}, \hat{W}'_{\ell/t})$ is a deterministic function of (Y, Y') , thus independent of (X, X') . Now we can use the induction hypothesis to conclude that the statement holds for the case of $a + 1$. Note that the total error is $O(\frac{\ell}{t}\epsilon) + \ell/t \cdot O(t\epsilon) = O(\ell\epsilon)$ since the part of $O(\frac{\ell}{t}\epsilon)$ decreases as a geometric sequence. Finally, the entropy requirement of Y is $(2^{\ell/t} - 1) \cdot 2m = (2^{\ell/t} - 1) \cdot 2 \cdot 2^{O(a \cdot t^{\frac{1}{a}})} \cdot s = 2^{\ell/t + O(a \cdot t^{\frac{1}{a}}) + 1} \cdot s$.

We now just need to choose a t to minimize this quantity. We can choose $t = \ell^{\frac{a}{a+1}}$ so that the entropy requirement of Y is $2^{O((a+1) \cdot \ell^{\frac{1}{a+1}})} \cdot s$. \square

We now have the following construction.

Construction 4.8. NIPM_x (which is strong in Y) or NIPM_y (which is strong in X).

Inputs:

- An error parameter $\epsilon > 0$ and a constant $a \in \mathbb{N}$.
- A random variable V supported on a boolean $L \times m$ matrix.
- A uniform string X on d_1 bits.
- A uniform string Y on d_2 bits.
- Let $d = c \log(\max\{d_1, d_2\}/\epsilon)$ for some constant $c > 1$.

Output: NIPM_x outputs a random variable $\mathbf{W}_x \in \{0, 1\}^m$, and NIPM_y outputs $\mathbf{W}_y \in \{0, 1\}^d$.

1. Let $\ell = \log L$.⁴ Let X_0 be a slice of X with length $4a \cdot d$, and Y_0 be a slice of Y with length $4a \cdot d$. Use X_0 and Y_0 to run an alternating extraction protocol, and output $(R_0, \dots, R_a) = \text{laExt}(X_0, Y_0)$ where each R_i has d bits.
2. Compute $Z = \text{Ext}(Y, R_0)$ and output $d_2/2$ bits, where Ext is the strong seeded extractor from Theorem 2.7.

⁴Without loss of generality we assume that L is a power of 2. Otherwise add 0 to the string until the length is a power of 2.

3. For every $i \in [L]$, compute $\overline{V}_i = \text{Ext}(Y_0, V_i)$ and output d bits. Then, compute $\hat{V}_i = \text{Ext}(X, \overline{V}_i)$ and output m bits.
4. Compute $(Z_1, \dots, Z_\ell) = \text{Conv}_a(Z, R_1, \dots, R_a)$ where each Z_i has d bits.
5. NIPM_x outputs $\mathbf{W}_x = \text{NIPM}(\hat{V}, Z_1, \dots, Z_\ell)$, where NIPM is the merger in Construction 4.4 and Lemma 4.5. NIPM_y outputs $\mathbf{W}_y = \text{Ext}(Y, \mathbf{W}_x)$ with d bits.

We now have the following lemma.

Lemma 4.9. *There exist a constant $c > 1$ such that for any $\epsilon > 0$ and any $L, m, d_1, d_2, n \in \mathbb{N}$ such that $d \geq c(\log \max\{d_1, d_2\} + \log(1/\epsilon))$, $m \geq d$, $d_1 \geq 8a \cdot d + 6m$ and $d_2 \geq 8a \cdot d + c^{a \cdot \log^{\frac{1}{a}} L} \cdot d$, the above construction gives an $(L, d_1, d_2, O(L\epsilon))$ -NIPM that is either strong in X or strong in Y .*

Proof. Note that Y_0 has min-entropy $4ad \geq 4d$, thus by Theorem 2.7 we have that for every $i \in [L]$,

$$(\overline{V}_i, V_i) \approx_\epsilon (U_d, V_i),$$

and there exists an $h \in [L]$ such that

$$(\overline{V}_h, \overline{V}'_h, V_h, V'_h) \approx_\epsilon (U_d, \overline{V}'_h, V_h, V'_h).$$

Note that conditioned on the fixing of (V, V') , we have that (X, X') and (Y, Y') are still independent, and furthermore $(\overline{V}, \overline{V}')$ is a deterministic function of (Y, Y') . Note that conditioned on the fixing of (X_0, X'_0) , the average conditional min-entropy of X is at least $8a \cdot d + 6m - 2 \cdot 4a \cdot d = 6m$. Thus again by Theorem 2.7 we have that for every $i \in [L]$,

$$(\hat{V}_i, \overline{V}_i) \approx_\epsilon (U_d, \overline{V}_i),$$

and there exists an $h \in [L]$ such that

$$(\hat{V}_h, \hat{V}'_h, \overline{V}_h, \overline{V}'_h) \approx_\epsilon (U_d, \hat{V}'_h, \overline{V}_h, \overline{V}'_h).$$

Note that now conditioned on the fixing of $(\overline{V}_h, \overline{V}'_h)$, we have that (X, X') and (Y, Y') are still independent, and furthermore (\hat{V}_h, \hat{V}'_h) is a deterministic function of (X, X') . Thus we basically have that conditioned on the fixing of (X_0, X'_0, Y_0, Y'_0) , (\hat{V}, \hat{V}') is a deterministic function of (X, X') and they satisfy the property needed by an NIPM.

Now, by Lemma 3.2, we have that

$$(Y_0, Y'_0, R_0, R'_0, \dots, R_a, R'_a) \approx_{O(a^2\epsilon)} (Y_0, Y'_0, U_d, R'_0, \dots, U_d, R'_a).$$

Note that conditioned on the fixing of (Y_0, Y'_0) , we have that (X, X') and (Y, Y') are still independent, and furthermore $(R_0, R'_0, \dots, R_a, R'_a)$ is a deterministic function of (X, X') . Also the average conditional min-entropy of Y is at least $d_2 - 2 \cdot 4a \cdot d = c^{a \cdot \log^{\frac{1}{a}} L} \cdot d > 3d_2/4$ for a large enough constant c . Thus by Theorem 2.7 we have that

$$(Z, R_0) \approx_\epsilon (U_{d_2/2}, R_0).$$

We can now fix (R_0, R'_0) . Note that now (Z_0, Z'_0) is a deterministic function of (Y, Y') , and $d_2/2 > \frac{1}{2}c^{a \cdot \log^{\frac{1}{a}} L} \cdot d$. Note that now $(R_1, R'_1, \dots, R_a, R'_a)$ still satisfies the look-ahead property. Thus as long as c is large enough, by Lemma 4.7 we have that

$$(Z_1, Z'_1, \dots, Z_\ell, Z'_\ell, X_0, X'_0) \approx_{O(\ell\epsilon)} (U_d, W'_1, \dots, U_d, W'_\ell, X_0, X'_0).$$

We can now fix (X_0, X'_0) , and note that conditioned on this fixing $(Z_1, Z'_1, \dots, Z_\ell, Z'_\ell)$ is a deterministic function of (Y, Y') . In summary, conditioned on the fixing of (X_0, X'_0, Y_0, Y'_0) , we have that (\hat{V}, \hat{V}') and $(Z_1, Z'_1, \dots, Z_\ell, Z'_\ell)$ satisfy the conditions required by Lemma 4.5. Therefore we can now apply that lemma to finish the proof. The total error is at most $O(L\epsilon) + O(a^2\epsilon) + O(\epsilon) + O(\ell\epsilon) = O(L\epsilon)$. \square

The extreme case of the above construction gives the following NIPM.

Construction 4.10. NIPM $_x$ (which is strong in Y) or NIPM $_y$ (which is strong in X).

Inputs:

- An error parameter $\epsilon > 0$.
- A random variable V supported on a boolean $L \times m$ matrix.
- A uniform string \mathbf{X} on n bits.
- A uniform string \mathbf{Y} on n' bits.

Output: NIPM $_x$ outputs a random variable $\mathbf{W}_x \in \{0, 1\}^m$, and NIPM $_y$ outputs $\mathbf{W}_y \in \{0, 1\}^{O(\log(n/\epsilon))}$.

1. Let $d_1 = c \log(n'/\epsilon)$ and $d_2 = c \log(n/\epsilon)$. Take a slice \mathbf{X}_0 of \mathbf{X} with length $10 \log \log L \cdot d_1$, and a slice \mathbf{Y}_0 of \mathbf{Y} with length $10 \log \log L \cdot d_2$.
2. Use \mathbf{X}_0 and \mathbf{Y}_0 to do an alternating extraction protocol, and output $(R_0, R_1, \dots, R_t) = \text{laExt}(\mathbf{X}_0, \mathbf{Y}_0)$ where $t = \log \log L$ and each R_i has $4d_1$ bits, each S_i (used in the alternating extraction) has d_2 bits.
3. For each $i \in [L]$, compute $\bar{\mathbf{Y}}_i = \text{Ext}(\mathbf{Y}_0, V_i)$ where each $\bar{\mathbf{Y}}_i$ outputs d_2 bits. Then compute $\bar{V}_i = \text{Ext}(\mathbf{X}, \bar{\mathbf{Y}}_i)$ where each \bar{V}_i outputs m bits. Here Ext is the strong seeded extractor from Theorem 2.7. Let \bar{V} be the matrix whose i 'th row is \bar{V}_i .
4. Let $\mathbf{Y}_1^0 = \mathbf{Y}$. For $j = 0$ to $\log \log L$ do the following. For $h = 1$ to 2^j , use \mathbf{Y}_h^j and R_j to do an alternating extraction protocol, and output $(S_{h1}^j, S_{h2}^j) = \text{laExt}(\mathbf{Y}_h^j, R_j)$, where each S_{hi}^j has $(\frac{\log^{\log a} L}{a^{j-1}} - 1)d_2$ bits. Note that altogether we get 2^{j+1} outputs and relabel them as $\mathbf{Y}_1^{j+1}, \dots, \mathbf{Y}_{2^{j+1}}^{j+1}$.
5. After the previous step, we get $2 \log L$ outputs. Let them be $\mathbf{Y}_1, \dots, \mathbf{Y}_{2 \log L}$, and output $\mathbf{W}_x = \text{NIPM}(\bar{V}, \mathbf{X}, \mathbf{Y}_1, \dots, \mathbf{Y}_{2 \log L})$ with m bits. Let $\mathbf{W}_y = \text{Ext}(\mathbf{Y}, \mathbf{W}_x)$ with d_2 bits.

We now have the following lemma.

Lemma 4.11. *There is a constant $c > 1$ such that suppose we have the following random variables and conditions:*

- V, V' , each supported on a boolean $L \times m$ matrix s.t for any $i \in [L]$, $V_i = U_m$. In addition, there exists an $h \in [L]$ such that $(V_h, V'_h) = (U_m, V'_h)$.
- \mathbf{Y}, \mathbf{Y}' , each supported on n' bits, where \mathbf{Y} is uniform.
- \mathbf{X}, \mathbf{X}' , each supported on n bits, where \mathbf{X} is uniform. In addition, \mathbf{X} is independent of (V, V') , and $(V, V', \mathbf{X}, \mathbf{X}')$ is independent of $(\mathbf{Y}, \mathbf{Y}')$.
- $m \geq c \log(n'/\epsilon)$, $n \geq 20c \log \log L \log(n'/\epsilon) + 6m$ and $n' \geq 20c \log^{\log a} L \log(n/\epsilon)$.

Let $(\mathbf{W}_x, \mathbf{W}_y)$ be the outputs of $(\text{NIPM}_x, \text{NIPM}_y)$ on $(V, \mathbf{X}, \mathbf{Y})$ and $(\mathbf{W}'_x, \mathbf{W}'_y)$ be the outputs of the $(\text{NIPM}_x, \text{NIPM}_y)$ on $(V', \mathbf{X}', \mathbf{Y}')$. Then

$$(\mathbf{W}_x, \mathbf{W}'_x, \mathbf{Y}, \mathbf{Y}') \approx_{O(L\epsilon)} (U_m, \mathbf{W}'_x, \mathbf{Y}, \mathbf{Y}')$$

and

$$(\mathbf{W}_y, \mathbf{W}'_y, V, V', \mathbf{X}, \mathbf{X}') \approx_{O(L\epsilon)} (U_{O(\log(n/\epsilon))}, \mathbf{W}'_y, V, V', \mathbf{X}, \mathbf{X}').$$

Proof. First, since $(V, V', \mathbf{X}, \mathbf{X}')$ is independent of $(\mathbf{Y}, \mathbf{Y}')$, as long as c is large enough, by Theorem 2.7 we know that for any $i \in [L]$,

$$(\bar{\mathbf{Y}}_i, V) \approx_\epsilon (U_d, V).$$

In addition, suppose for some $h \in [L]$ we have that $(V_h, V'_h) = (U_m, V'_h)$, then we can first fix V'_h and then $\bar{\mathbf{Y}}_h$. Conditioned on this fixing V_h is still uniform, the average conditional min-entropy of \mathbf{Y}_0 is at least $10 \log \log L \cdot d - d > 3d$ and V_h and \mathbf{Y}_0 are still independent, thus by Theorem 2.7 we have that

$$(\bar{\mathbf{Y}}_h, \bar{\mathbf{Y}}'_h, V, V') \approx_\epsilon (U_d, \bar{\mathbf{Y}}'_h, V, V').$$

In other words, the random variables $\{(\bar{\mathbf{Y}}_i, \bar{\mathbf{Y}}'_i)\}$ inherit the properties of $\{(V_i, V'_i)\}$. We now ignore the errors since this adds at most $L\epsilon$ to the final error. Now we fix (V, V') . Note that conditioned on this fixing, the random variables $(\bar{\mathbf{Y}}_i, \bar{\mathbf{Y}}'_i)$ are deterministic functions of $(\mathbf{Y}_0, \mathbf{Y}'_0)$, and are thus independent of $(\mathbf{X}, \mathbf{X}')$. Furthermore, we have that conditioned on this fixing, \mathbf{X} is still uniform. In addition, even conditioned on the fixing of $(\mathbf{X}_0, \mathbf{X}'_0)$, the average conditional min-entropy of \mathbf{X} is at least $20c \log \log L \log(n'/\epsilon) + 6m - 2 \cdot 10 \log \log L \cdot d_1 = 6m$. Thus by the same argument before we have that for any $i \in [L]$,

$$(\bar{V}_i, \mathbf{Y}_0, \mathbf{X}_0, \mathbf{X}'_0) \approx_\epsilon (U_m, \mathbf{Y}_0, \mathbf{X}_0, \mathbf{X}'_0),$$

and that there exists an $h \in [L]$ such that

$$(\bar{V}_h, \bar{V}'_h, \mathbf{Y}_0, \mathbf{Y}'_0, \mathbf{X}_0, \mathbf{X}'_0) \approx_\epsilon (U_m, \bar{V}'_h, \mathbf{Y}_0, \mathbf{Y}'_0, \mathbf{X}_0, \mathbf{X}'_0).$$

We will again ignore the error for now since this adds at most $L\epsilon$ to the final error. Next, by Lemma 3.2 we have that for any $0 \leq j \leq t-1$,

$$(R_{j+1}, (R_1, R'_1, \dots, R_j, R'_j), \mathbf{Y}_0, \mathbf{Y}'_0) \approx_{O(t\epsilon)} (U_{4d_1}, (R_1, R'_1, \dots, R_j, R'_j), \mathbf{Y}_0, \mathbf{Y}'_0).$$

Thus by a hybrid argument and the triangle inequality, we have that

$$(\mathbf{Y}_0, \mathbf{Y}'_0, R_1, R'_1, \dots, R_t, R'_t) \approx_{O(t^2\epsilon)} (\mathbf{Y}_0, \mathbf{Y}'_0, U_{4d_1}, R'_1, \dots, U_{4d_1}, R'_t),$$

where each U_{4d_1} is independent of all the previous random variables (but may depend on later random variables). From now on, we will proceed as if each R_j is uniform given $(\mathbf{Y}_0, \mathbf{Y}'_0, \{R_1, R'_1, \dots, R_{j-1}, R'_{j-1}\})$, since this only adds $O(t^2\epsilon)$ to the final error.

Now we can fix $(\mathbf{Y}_0, \mathbf{Y}'_0)$. Note that conditioned on this fixing, $(\bar{V}, \bar{V}', R_1, R'_1, \dots, R_t, R'_t)$ are deterministic functions of $(V, V', \mathbf{X}, \mathbf{X}')$, and thus independent of $(\mathbf{Y}, \mathbf{Y}')$. Also note that conditioned on this fixing, the average conditional min-entropy of \mathbf{Y} is at least $20 \log^{\log a} L \cdot d_2 - 2 \cdot 10 \log \log L \cdot d_2 > a^2 \log^{\log a} L \cdot d_2$. We now prove the following claim.

Claim 4.12. *Let $\bar{R}_j = (R_1, \dots, R_j)$. Suppose that at the beginning of the j 'th iteration, we have that conditioned on the fixing of \bar{R}_{j-1} , the following holds.*

1. $(\mathbf{X}, \mathbf{X}')$ is independent of $(\mathbf{Y}, \mathbf{Y}')$, and $(\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{2^j}, \mathbf{Y}'_{2^j})$ is a deterministic function of $(\mathbf{Y}, \mathbf{Y}')$.
2. For every $h \in [2^j]$, the average conditional min-entropy of \mathbf{Y}_h given $(\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{h-1}, \mathbf{Y}'_{h-1})$ is at least $(\frac{\log^{\log a} L}{a^{j-2}} - 1)d_2$.

Then at the end of the j 'th iteration, the following holds.

1. Conditioned on the fixing of \bar{R}_j , $(\mathbf{X}, \mathbf{X}')$ is independent of $(\mathbf{Y}, \mathbf{Y}')$, and $(\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{2^{j+1}}, \mathbf{Y}'_{2^{j+1}})$ is a deterministic function of $(\mathbf{Y}, \mathbf{Y}')$.
2. For every $h \in [2^{j+1}]$,

$$(\mathbf{Y}_h, (\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{h-1}, \mathbf{Y}'_{h-1}), \bar{R}_j) \approx_\epsilon (U_{(\frac{\log^{\log a} L}{a^{j-1}} - 1)d_2}, (\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{h-1}, \mathbf{Y}'_{h-1}), \bar{R}_j).$$

Proof of the claim. First, since the computation in the j 'th iteration only involves (R_j, R'_j) and $(\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{2^j}, \mathbf{Y}'_{2^j})$, and (R_j, R'_j) is a deterministic function of $(\mathbf{X}, \mathbf{X}')$ conditioned on the fixing of the previous random variables, we know that at the end of the j 'th iteration, conditioned on the fixing of (R_1, \dots, R_j) we have that $(\mathbf{X}, \mathbf{X}')$ is independent of $(\mathbf{Y}, \mathbf{Y}')$, and $(\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{2^{j+1}}, \mathbf{Y}'_{2^{j+1}})$ is a deterministic function of $(\mathbf{Y}, \mathbf{Y}')$.

Next, we use $(Z_1, Z'_1, \dots, Z_{2^{j+1}}, Z'_{2^{j+1}})$ to represent the outputs computed from (R_j, R'_j) and $(\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{2^j}, \mathbf{Y}'_{2^j})$, and assume that $2\ell - 1 \leq h \leq 2\ell$ for some ℓ , then Z_h is obtained from \mathbf{Y}_ℓ . We can now first fix $(\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{\ell-1}, \mathbf{Y}'_{\ell-1})$, and conditioned on this fixing \mathbf{Y}_ℓ has average conditional min-entropy at least $(\frac{\log^{\log a} L}{a^{j-2}} - 1)d_2$. Now by Lemma 3.2 we have that

$$(S_1^\ell, R_j, R'_j) \approx_\epsilon (U_{(\frac{\log^{\log a} L}{a^{j-1}} - 1)d_2}, R_j, R'_j)$$

and

$$(S_2^\ell, S_1^\ell, S_1^{\ell'}, R_j, R'_j) \approx_\epsilon (U_{(\frac{\log^{\log a} L}{a^{j-1}} - 1)d_2}, S_1^\ell, S_1^{\ell'}, R_j, R'_j),$$

since $(\frac{\log^{\log a} L}{a^{j-2}} - 1)d_2 \geq 2 \cdot (\frac{\log^{\log a} L}{a^{j-1}} - 1)d_2 + (1 + \alpha)(\frac{\log^{\log a} L}{a^{j-1}} - 1)d_2 + d_2$ and $4d_1 \geq 2d_1 + 1.1d_1 + 0.9d_1$. Thus as long as the constant c is large enough one can make sure that $\min\{d_2, 0.9d_1\} \geq 2 \log(1/\epsilon)$, and we can extract $(\frac{\log^{\log a} L}{a^{j-1}} - 1)d_2$ bits from entropy $(1 + \alpha)(\frac{\log^{\log a} L}{a^{j-1}} - 1)d_2$ and d_1 bits from entropy $1.1d_1$. Note that $(Z_1, Z'_1, \dots, Z_{2\ell-2}, Z'_{2\ell-2})$ are computed from $(\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{\ell-1}, \mathbf{Y}'_{\ell-1})$ and (R_j, R'_j) , and $(\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{\ell-1}, \mathbf{Y}'_{\ell-1})$ are already fixed. Thus the second part of the claim also holds. \square

Now note that at the beginning of the first iteration, the condition of the claim holds. Thus if we ignore the errors, then we can apply the claim repeatedly until the end of the iteration. At this time for each $h \in [\log L]$ we have that \mathbf{Y}_h has at least $(\frac{\log^{\log a} L}{a^{\log L - 1}} - 1)d_2 > d_2$ bits. Furthermore

$$(\mathbf{Y}_h, (\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{h-1}, \mathbf{Y}'_{h-1}), \overline{R}_t) \approx (U, (\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{h-1}, \mathbf{Y}'_{h-1}), \overline{R}_t).$$

The total error so far is $O(L\epsilon) + O(t^2\epsilon) + \sum_{j=0}^{\log \log L} 2^j \cdot 2\epsilon = O(L\epsilon)$. Note that now conditioned on all the fixed random variables $(\mathbf{X}_0, \mathbf{X}'_0, \mathbf{Y}_0, \mathbf{Y}'_0, \overline{R}_t)$ (note that \overline{R}_t is a deterministic function of $(\mathbf{X}_0, \mathbf{X}'_0, \mathbf{Y}_0, \mathbf{Y}'_0)$), we have that $(V, V', \mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{2\log L}, \mathbf{Y}'_{2\log L}, \mathbf{X}, \mathbf{X}')$ satisfies the conditions of the Lemma 4.5, since the average conditional min-entropy of X is at least $n - 20 \log \log L \cdot d_1 \geq 6m$. Now we can apply Lemma 4.5 to show that

$$(\mathbf{W}_x, \mathbf{W}'_x, \mathbf{Y}, \mathbf{Y}') \approx (U_m, \mathbf{W}'_x, \mathbf{Y}, \mathbf{Y}'),$$

where the total error is $O(L\epsilon) + O(L\epsilon) = O(L\epsilon)$. Furthermore, note that conditioned on the fixing of $(\mathbf{Y}_1, \mathbf{Y}'_1, \dots, \mathbf{Y}_{2\log L}, \mathbf{Y}'_{2\log L})$, we have that $(\mathbf{W}_x, \mathbf{W}'_x)$ is a deterministic function of $(V, V', \mathbf{X}, \mathbf{X}')$, and thus independent of $(\mathbf{Y}, \mathbf{Y}')$. Also note that \mathbf{Y} has average conditional min-entropy at least $20c \log^{\log a} L \log(n/\epsilon) - 4 \log L d_2 > 10d_2$. Thus by Theorem 2.7 we have that

$$(\mathbf{W}_y, \mathbf{W}'_y, \mathbf{W}_x, \mathbf{W}'_x) \approx (U_{d_2}, \mathbf{W}'_y, \mathbf{W}_x, \mathbf{W}'_x),$$

where the error is $O(L\epsilon) + O(\epsilon) = O(L\epsilon)$. Note that given $(\mathbf{W}_x, \mathbf{W}'_x)$, we have that $(\mathbf{W}_y, \mathbf{W}'_y)$ is a deterministic function of $(\mathbf{Y}, \mathbf{Y}')$. Thus we also have that

$$(\mathbf{W}_y, \mathbf{W}'_y, V, V', \mathbf{X}, \mathbf{X}') \approx_{O(L\epsilon)} (U_{d_2}, \mathbf{W}'_y, V, V', \mathbf{X}, \mathbf{X}').$$

\square

5 Correlation Breaker with Advice

We now use our non-malleable independence preserving mergers to construct improved correlation breakers with advice. A correlation breaker uses independent randomness to break the correlations between several correlated random variables. The first correlation breaker appears implicitly in the author's work [Lil13a], and this object is strengthened and formally defined in [Coh15]. A correlation breaker with advice additionally uses some string as an advice. This object was first introduced and used without its name in [CGL16], and then explicitly defined in [Coh16b].

Definition 5.1 (Correlation breaker with advice). A function

$$\text{AdvCB} : \{0, 1\}^n \times \{0, 1\}^d \times \{0, 1\}^L \rightarrow \{0, 1\}^m$$

is called a (k, k', ε) -correlation breaker with advice if the following holds. Let Y, Y' be d -bit random variables such that $H_\infty(Y) \geq k'$. Let X, X' be n -bit random variables with $H_\infty(X) \geq k$, such that (X, X') is independent of (Y, Y') . Then, for any pair of distinct L -bit strings α, α' ,

$$(\text{AdvCB}(X, Y, \alpha), \text{AdvCB}(X', Y', \alpha')) \approx_\varepsilon (U, \text{AdvCB}(X', Y', \alpha')).$$

In addition, we say that AdvCB is strong if

$$\begin{aligned} & (\text{AdvCB}(X, Y, \alpha), \text{AdvCB}(X', Y', \alpha'), Y, Y') \\ & \approx_\varepsilon (U, \text{AdvCB}(X', Y', \alpha'), Y, Y'). \end{aligned}$$

Our construction needs the following flip-flop extraction scheme, which was constructed by Cohen [Coh15] using alternating extraction, based on a previous similar construction of the author [Li13a]. The flip-flop function can be viewed as a basic correlation breaker, which (informally) uses an independent source \mathbf{X} to break the correlation between two r.v's \mathbf{Y} and \mathbf{Y}' , given an advice bit.

Theorem 5.2 ([Coh15, CGL16]). *There exists a constant $c_{5.2}$ such that for all $n > 0$ and any $\varepsilon > 0$, there exists an explicit function $\text{flip-flop} : \{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^m$, $m = 0.4k$, satisfying the following: Let \mathbf{X} be an (n, k) -source, and \mathbf{X}' be a random variable on n bits arbitrarily correlated with \mathbf{X} . Let \mathbf{Y} be an independent uniform seed on d bits, and \mathbf{Y}' be a random variable on d bits arbitrarily correlated with \mathbf{Y} . Suppose $(\mathbf{X}, \mathbf{X}')$ is independent of $(\mathbf{Y}, \mathbf{Y}')$. If $k, d \geq C_{5.2} \log(n/\varepsilon)$, then for any bit b ,*

$$|\text{flip-flop}(\mathbf{X}, \mathbf{Y}, b), \mathbf{Y}, \mathbf{Y}' - U_m, \mathbf{Y}, \mathbf{Y}'| \leq \varepsilon.$$

Furthermore, for any bits b, b' with $b \neq b'$,

$$\begin{aligned} & |\text{flip-flop}(\mathbf{X}, \mathbf{Y}, b), \text{flip-flop}(\mathbf{X}', \mathbf{Y}', b'), \mathbf{Y}, \mathbf{Y}' \\ & - U_m, \text{flip-flop}(\mathbf{X}', \mathbf{Y}', b'), \mathbf{Y}, \mathbf{Y}'| \leq \varepsilon. \end{aligned}$$

5.1 Asymmetric correlation breaker

We will present correlation breakers that use general NIPMs. By plugging in various NIPMs this gives different correlation breakers.

Construction 5.3. Inputs:

- Let $\ell, m \in \mathbb{N}$ be two integers, $\varepsilon > 0$ be an error parameter.
 - X, Y , two independent sources on n bits and s bits respectively, with min-entropy at least $n - \ell$ and $s - \ell$.
 - an advice string $\alpha \in \{0, 1\}^L$.
 - An $(L, d_1, d_2, O(L\varepsilon))$ -NIPM $_x$ that is strong in Y .
 - Let IP be the two source extractor from Theorem 2.8.
1. Let $d' = O(\log(\max\{n, s\}/\varepsilon))$ be the seed length of the extractor from Theorem 2.7, and let $d = 8d'$. Let X^0 be a slice of X with length $d + 2\ell + 2\log(1/\varepsilon)$, and Y^0 be a slice of Y with length $d + 2\ell + 2\log(1/\varepsilon)$.

2. Compute $Z = \text{IP}(X^0, Y^0)$ and output d bits.
3. Use X and Z to do an alternating extraction, and output two random variables $(X_0, X_1) = \text{laExt}(X, Z)$ where each X_i has $3m$ bits.
4. Use Y and Z to do an alternating extraction, and output two random variables $(Y_0, Y_1) = \text{laExt}(Y, Z)$ where each Y_i has $3d$ bits.
5. Use X_1, Y_1, α to obtain an $L \times m$ matrix V , where for any $i \in [L]$, $V_i = \text{flip-flop}(X_1, Y_1, \alpha_i)$ and outputs m bits.
6. Compute $\hat{X} = \text{Ext}(X, Y_0)$ and output $n/2$ bits. Compute $\hat{Y} = \text{Ext}(Y, X_0)$ and output $s/2$ bits. Here Ext is the strong seeded extractor from Theorem 2.7.
7. Output $\hat{V} = \text{NIPM}_x(V, \hat{X}, \hat{Y})$.

We now have the following lemma.

Lemma 5.4. *There exists a constant $c > 1$ such that the following holds. Suppose that there exists an $(L, d_1, d_2, O(L\epsilon))$ -NIPM that is strong in Y which outputs m bits, then there exists an explicit $(n - \ell, s - \ell, O(L\epsilon))$ AdvCB : $\{0, 1\}^n \times \{0, 1\}^s \times \{0, 1\}^L \rightarrow \{0, 1\}^m$ as long as $m \geq c \log(\max\{n, s\}/\epsilon)$, $n \geq 20m + 2d_1 + 5\ell + 4 \log(1/\epsilon)$ and $s \geq m + 2d_2 + 5\ell + 4 \log(1/\epsilon)$.*

Proof. Throughout the proof we will use letters with prime to denote the corresponding random variables obtained from (X', Y', α') . First, notice that both X^0 and Y^0 have min-entropy at least $d + \ell + 2 \log(1/\epsilon)$. Thus by Theorem 2.8 we have that

$$(Z, X^0) \approx_\epsilon (U_d, X^0)$$

and

$$(Z, Y^0) \approx_\epsilon (U_d, Y^0).$$

We now ignore the error ϵ . Note that conditioned on the fixing of (X^0, X'^0) , (Z, Z') is a deterministic function of (Y^0, Y'^0) , and thus independent of (X, X') . Moreover, the average conditional min-entropy of X given this fixing is at least $n - \ell - 2(d + 2\ell + 2 \log(1/\epsilon)) \geq 10m$ as long as c is large enough. Thus by Lemma 3.2 (note that the extractor from Z side can use seed length d') we have that

$$(Y^0, Y'^0, X_0, X'_0, X_1, X'_1, Z, Z') \approx_{O(\epsilon)} (Y^0, Y'^0, U_{3m}, X'_0, U_{d_1}, X'_1, Z, Z'),$$

where each U_{3m} is uniform given the previous random variables, but may depend on later random variables. Similarly, note that conditioned on the fixing of (Y^0, Y'^0) , (Z, Z') is a deterministic function of (X^0, X'^0) , and thus independent of (Y, Y') . Moreover, the average conditional min-entropy of Y given this fixing is at least $s - \ell - 2(d + 2\ell + 2 \log(1/\epsilon)) \geq 10d$. Thus by Lemma 3.2 we have that

$$(Y_0, Y'_0, Y_1, Y'_1, Z, Z', X^0, X'^0) \approx_{O(\epsilon)} (U_{3d}, Y'_0, U_{d_2}, Y'_1, Z, Z', X^0, X'^0),$$

where each U_{3d} is uniform given the previous random variables, but may depend on later random variables. We can now fix (X^0, X'^0, Y^0, Y'^0) , and conditioned on this fixing, we have that

(X, X') and (Y, Y') are still independent, (X_0, X'_0, X_1, X'_1) is a deterministic function of (X, X') , and (Y_0, Y'_0, Y_1, Y'_1) is a deterministic function of (Y, Y') . Further they satisfy the look-ahead properties in the previous two equations. We will ignore the error for now since this only adds at most $O(\epsilon)$ to the final error.

We now claim that conditioned on the fixing of $(X_0, X'_0, Y_0, Y'_0, Y_1, Y'_1)$ (and ignoring the error), the random variables $(V, V', \hat{X}, \hat{X}')$ and (\hat{Y}, \hat{Y}') satisfy the conditions required by Lemma 4.9. To see this, note that if we fix (Y_0, Y'_0, Y_1, Y'_1) , then the average conditional min-entropy of Y is at least $s - \ell - 2(d + 2\ell + 2\log(1/\epsilon)) - 2 \cdot 3d > 2s/3$ as long as c is large enough. Thus by Theorem 2.7 we have that

$$(\hat{Y}, X_0, X'_0) \approx_\epsilon (U_{s/2}, X_0, X'_0).$$

Thus conditioned on the further fixing of (X_0, X'_0) , we have that (\hat{Y}, \hat{Y}') is a deterministic function of (Y, Y') , and $s/2 \geq d_2$. On the other hand, conditioned on the fixing of (X_0, X'_0) and (Y_0, Y'_0) , we have X_1 is still close to uniform. Thus by Theorem 5.2 we have that for any $i \in [L]$,

$$|V_i, Y_1, Y'_1 - U_m, Y_1, Y'_1| \leq \epsilon$$

and there exists $i \in [L]$ such that

$$|V_i, V'_i, Y_1, Y'_1 - U_m, V'_i, Y_1, Y'_1| \leq \epsilon.$$

We now further fix (Y_1, Y'_1) . Note that conditioned on this fixing (X, X') and (Y, Y') are still independent. Furthermore (V, V') is now a deterministic function of (X_1, X'_1) , and thus independent of (Y, Y') . Finally, note that conditioned on the fixing of (X_0, X'_0, X_1, X'_1) , the average conditional min-entropy of X is at least $n - \ell - 2(d + 2\ell + 2\log(1/\epsilon)) - 2 \cdot 3m > 2n/3$. Thus by Theorem 2.7 we have that

$$(\hat{X}, Y_0, Y'_0) \approx_\epsilon (U_{n/2}, Y_0, Y'_0).$$

Thus conditioned on the further fixing of (Y_0, Y'_0) , we have that (\hat{X}, \hat{X}') is a deterministic function of (X, X') , and $n/2 \geq d_1$. Thus, even if conditioned on the fixing of $(X_0, X'_0, X_1, X'_1, Y_0, Y'_0, Y_1, Y'_1)$, we have that $(\hat{X}$ is close to $U_{n/2}$. Since (V, V') is obtained from (X_1, X'_1, Y_1, Y'_1) , we know that $(\hat{X}$ is close to uniform even given $(X_0, X'_0, Y_0, Y'_0, Y_1, Y'_1)$ and (V, V') . Thus by Lemma 4.9 we have that

$$(\hat{V}, \hat{V}', Y, Y') \approx (U_m, \hat{V}', Y, Y'),$$

where the error is $O(L\epsilon) + O(L\epsilon) + O(\epsilon) = O(L\epsilon)$. □

Next we give another correlation breaker, which recycles the randomness used.

Construction 5.5. Inputs:

- Let $\ell, m \in \mathbb{N}$ be two integers, $\epsilon > 0$ be an error parameter.
- X, Y , two independent sources on n bits with min-entropy at least $n - \ell$.
- an advice string $\alpha \in \{0, 1\}^L$ and an integer $2 \leq t \leq L$.
- An $(L, d_1, d_2, O(L\epsilon))$ -NIPM $_y$ that is strong in X .

- Let IP be the two source extractor from Theorem 2.8.
1. Let $d' = O(\log(n/\epsilon))$ be the seed length of the extractor from Theorem 2.7, and let $d = 8 \frac{\log L}{\log t} d'$. Let X^0 be a slice of X with length $d + 2\ell + 2\log(1/\epsilon)$, and Y^0 be a slice of Y with length $d + 2\ell + 2\log(1/\epsilon)$.
 2. Compute $Z = \text{IP}(X^0, Y^0)$ and output d bits.
 3. Use X and Z to do an alternating extraction, and output $3 \frac{\log L}{\log t} + 1$ random variables $X_0, \dots, X_{3 \frac{\log L}{\log t}}$ where each X_i has d_1 bits.
 4. Use Y and Z to do an alternating extraction, and output two random variables Y_0, Y_1 where each Y_i has d_2 bits.
 5. Use X_0, Y_0, α to obtain an $L \times m$ matrix V^0 , where for any $i \in [L]$, $V_i^0 = \text{flip-flop}(X_0, Y_0, \alpha_i)$ and outputs m bits.
 6. For $i = 1$ to $\frac{\log L}{\log t}$ do the following. Merge every t rows of V^{i-1} using NIPM_y and (X_{3i-2}, Y_i) , and output d' bits. Concatenate the outputs to become another matrix W^i . Note that W^i has L/t^i rows. Then for every row $j \in [L/t^i]$, compute $V_j^i = \text{Ext}(X_{3i}, W_j^i)$ to obtain a new matrix V^i . Finally let $Y_{i+1} = \text{Ext}(Y, X_{3i-1})$ and output d_2 bits.
 7. Output $\hat{V} = V^{\frac{\log L}{\log t}}$.

We now have the following lemma.

Lemma 5.6. *There exists a constant $c > 1$ such that the following holds. Suppose that for any $t \in \mathbb{N}$ there exists an $(t, d_1, d_2, O(t\epsilon))$ - NIPM_y that is strong in X which outputs $d' = O(\log(n/\epsilon))$ bits, then there exists an explicit $(n - \ell, n - \ell, O(L\epsilon))$ correlation breaker with advice $\text{AdvCB} : \{0, 1\}^n \times \{0, 1\}^n \times \{0, 1\}^L \rightarrow \{0, 1\}^m$ as long as $d_1 \geq 4m$, $m \geq c \log(d_2/\epsilon)$, and $n \geq c \frac{\log L}{\log t} \log(n/\epsilon) + \max\{8 \frac{\log L}{\log t} d_1, 2t \cdot d' + 4d_2\} + 5\ell + 4\log(1/\epsilon)$.*

Proof. Throughout the proof we will use letters with prime to denote the corresponding random variables obtained from (X', Y', α') . First, notice that both X^0 and Y^0 have min-entropy at least $d + \ell + 2\log(1/\epsilon)$. Thus by Theorem 2.8 we have that

$$(Z, X^0) \approx_\epsilon (U_d, X^0)$$

and

$$(Z, Y^0) \approx_\epsilon (U_d, Y^0).$$

We now ignore the error ϵ . Note that conditioned on the fixing of (X^0, X'^0) , (Z, Z') is a deterministic function of (Y^0, Y'^0) , and thus independent of (X, X') . Moreover, the average conditional min-entropy of X given this fixing is at least $n - \ell - 2(d + 2\ell + 2\log(1/\epsilon)) \geq 8 \frac{\log L}{\log t} d_1$ as long as c is large enough. Thus by Lemma 3.2 (note that the extractor from Z side can use seed length d') we have that

$$(Y^0, Y'^0, Z, Z', X_0, X'_0, \dots, X_{3 \frac{\log L}{\log t}}, X'_{3 \frac{\log L}{\log t}}) \approx_{O((\frac{\log L}{\log t})^2 \epsilon)} (Y^0, Y'^0, Z, Z', U_{d_1}, X'_0, \dots, U_{d_1}, X'_{3 \frac{\log L}{\log t}}),$$

where each U_{d_1} is uniform given the previous random variables, but may depend on later random variables. Similarly, note that conditioned on the fixing of (Y^0, Y'^0) , (Z, Z') is a deterministic function of (X^0, X'^0) , and thus independent of (Y, Y') . Moreover, the average conditional min-entropy of Y given this fixing is at least $n - \ell - 2(d + 2\ell + 2 \log(1/\epsilon)) \geq 4d_2$. Thus by Lemma 3.2 we have that

$$(Z, Z', X^0, X'^0, Y_0, Y'_0, Y_1, Y'_1) \approx_{O(\epsilon)} (Z, Z', X^0, X'^0, U_{d_2}, Y'_0, U_{d_2}),$$

where each U_{d_2} is uniform given the previous random variables, but may depend on later random variables. We can now fix (X^0, X'^0, Y^0, Y'^0) , and conditioned on this fixing, we have that (X, X') and (Y, Y') are still independent, $(X_0, X'_0, \dots, X_{3 \frac{\log L}{\log t}}, X'_{3 \frac{\log L}{\log t}})$ is a deterministic function of (X, X') , and (Y_0, Y'_0, Y_1, Y'_1) is a deterministic function of (Y, Y') . Further they satisfy the look-ahead properties in the previous two equations. We will ignore the error for now since this only adds at most $O((\frac{\log L}{\log t})^2 \epsilon)$ to the final error.

Now by Theorem 5.2 we have that for any $i \in [L]$,

$$|V_i^0, Y_0, Y'_0 - U_m, Y_0, Y'_0| \leq \epsilon$$

and there exists $i \in [L]$ such that

$$|V_i^0, V_i'^0, Y_0, Y'_0 - U_m, V_i'^0, Y_0, Y'_0| \leq \epsilon.$$

We now further fix (Y_0, Y'_0) . Note that conditioned on this fixing (X, X') and (Y, Y') are still independent. Furthermore (V^0, V'^0) is now a deterministic function of (X_0, X'_0) , and thus independent of (Y, Y') . Thus by the property of NIPM $_y$ we have that for every row j in W^1 ,

$$(W_j^1, V^0, V'^0, X_1, X'_1) \approx_{O(t\epsilon)} (U_{d'}, V^0, V'^0, X_1, X'_1),$$

and there exists a row j such that

$$(W_j^1, W_j'^1, V^0, V'^0, X_1, X'_1) \approx_{O(t\epsilon)} (U_{d'}, W_j'^1, V^0, V'^0, X_1, X'_1).$$

Note that we have fixed (X^0, X'^0, Y^0, Y'^0) , and if we further condition on the fixing of $(X_0, X'_0, Y_0, Y'_0, X_1, X'_1)$, then (W^1, W'^1) is a deterministic function of (Y, Y') . Furthermore (X, X') and (Y, Y') are still independent. We will now use induction to prove the following claim (note that we have already fixed (X^0, X'^0, Y^0, Y'^0)).

Claim 5.7. *Let $T_i = (Y_0, Y'_0, X_0, X'_0, \dots, X_{3i-2}, X'_{3i-2})$. In the i 'th iteration, the following holds.*

1. *Conditioned on the further fixing of T_i , we have that (X, X') and (Y, Y') are still independent, and furthermore (W^i, W'^i) is a deterministic function of (Y, Y') .*
2. *For every row j in W^i ,*

$$(W_j^i, T_i) \approx_{\epsilon_i} (U_{d'}, T_i),$$

and there exists a row j such that

$$(W_j^i, W_j'^i, T_i) \approx_{\epsilon_i} (U_{d'}, W_j'^i, T_i),$$

where $\epsilon_i = O(\sum_{j=1}^i t^j \epsilon)$.

Proof of the claim. The base case of $i = 1$ is already proved above. Now suppose the claim holds for the i 'th iteration, we show that it also holds for the $i + 1$ 'th iteration.

To see this, note that conditioned on the fixing of T_i , (X, X') and (Y, Y') are still independent, and furthermore (W^i, W'^i) is a deterministic function of (Y, Y') and thus independent of (X, X') . Note that Y_{i+1} is computed from Y and X_{3i-1} while V^i is computed from X_{3i} and W^i . Thus if we further fix X_{3i-1}, X'_{3i-1} and (W^i, W'^i) , then (X, X') and (Y, Y') are still independent, and furthermore Y_{i+1} is a deterministic function of Y and V^i is a deterministic function of X_{3i} . Now W^{i+1} is computed from V^i, X_{3i+1} and Y_{i+1} . Thus if we further fix (X_{3i}, X'_{3i}) and (X_{3i+1}, X'_{3i+1}) (i.e., we have fixed T_{i+1}) then (X, X') and (Y, Y') are still independent, and furthermore (W^{i+1}, W'^{i+1}) is a deterministic function of (Y, Y') .

Next, let h be the row in W^i such that

$$(W_h^i, W_h'^i, T_i) \approx_{\epsilon_i} (U_{d'}, W_h'^i, T_i).$$

Note that V^i has the same number of rows as W^i , and consider the merging of some t rows in V^i that contain row h into W_j^{i+1} (the merging of the other rows is similar and simpler). Without loss of generality assume that these t rows are row $1, 2, \dots, t$.

First, since for every row j in W^i ,

$$(W_j^i, T_i) \approx_{\epsilon_i} (U_{d'}, T_i),$$

and rows h in W^i and W'^i satisfy the independence property, by Theorem 2.7 (and ignoring the error ϵ_i) we have that for every $j \in [t]$,

$$(V_j^i, T_i, X_{3i-1}, X'_{3i-1}, W_j^i, W_j'^i) \approx_{\epsilon} (U_m, T_i, X_{3i-1}, X'_{3i-1}, W_j^i, W_j'^i),$$

and

$$(V_h^i, V_h'^i, T_i, X_{3i-1}, X'_{3i-1}, W_j^i, W_j'^i) \approx_{\epsilon} (U_m, V_h^i, T_i, X_{3i-1}, X'_{3i-1}, W_j^i, W_j'^i).$$

This is because X_{3i} has average conditional min-entropy at least d_1 even conditioned on the fixing of (X_{3i-1}, X'_{3i-1}) . We now ignore the error ϵ . Note that conditioned on the fixing of $(W_j^i, W_j'^i)$, we have that $(V_j^i, V_j'^i)$ is a deterministic function of (X_{3i}, X'_{3i}) , and thus independent of (Y, Y') . We now fix $\{(W_j^i, W_j'^i), j \in [t]\}$. Note that conditioned on this fixing $\{V_j^i, j \in [t]\}$ and $\{V_j'^i, j \in [t]\}$ each is a $t \times m$ matrix, and a deterministic function of (X_{3i}, X'_{3i}) . Further note that they form two matrices that meet the condition to apply an NIPM. Since $\{(W_j^i, W_j'^i), j \in [t]\}$ is a deterministic function of (Y, Y') , conditioned on this fixing (X, X') and (Y, Y') are still independent. Furthermore the average conditional min-entropy of Y is at least $n - \ell - 2(d + 2\ell + 2 \log(1/\epsilon)) - 2d_2 - 2td' \geq 2d_2$. Thus by Theorem 2.7 we have that

$$(Y_{i+1}, X_{3i-1}) \approx_{\epsilon} (U_{d_2}, X_{3i-1}).$$

Note that conditioned on the fixing of X_{3i-1} , we have that Y_{i+1} is a deterministic function of Y . Thus we can now further fix (X_{3i-1}, X'_{3i-1}) , and conditioned on this fixing, Y_{i+1} is still close to uniform. To conclude, now conditioned on the fixing of $\{(W_j^i, W_j'^i), j \in [t]\}$ and (X_{3i-1}, X'_{3i-1}) , we have that $\{V_j^i, j \in [t]\}$ and $\{V_j'^i, j \in [t]\}$ each is a $t \times m$ matrix, and a deterministic function of (X_{3i}, X'_{3i}) ; Y_{i+1} is still close to uniform and (Y_{i+1}, Y'_{i+1}) is a deterministic function of (Y, Y') . Furthermore X_{3i+1} is close to uniform. Now we can use the property of NIPM_y to show that after merging these t rows, the corresponding row j in W^{i+1} satisfies

$$\begin{aligned} & (W_j^{i+1}, W_j'^{i+1}, T_i, X_{3i-1}, X'_{3i-1}, X_{3i}, X'_{3i}, X_{3i+1}, X'_{3i+1}) \\ & \approx_{t\epsilon} (U_d, W_j^{i+1}, T_i, X_{3i-1}, X'_{3i-1}, X_{3i}, X'_{3i}, X_{3i+1}, X'_{3i+1}). \end{aligned}$$

Adding back all the errors we get that

$$(W_j^{i+1}, W_j'^{i+1}, T_{i+1}) \approx_{\epsilon_{i+1}} (U_d, W_j^{i+1}, T_{i+1}),$$

where $\epsilon_{i+1} = t\epsilon_i + O(t\epsilon) = O(\sum_{j=1}^{i+1} t^j \epsilon)$. □

Now we are basically done. In the last iteration we know that $W^{\frac{\log L}{\log t}}$ has reduced to one row, and $W^{\frac{\log L}{\log t}}$ is close to uniform given $W'^{\frac{\log L}{\log t}}$. Also conditioned on the fixing of $T^{\frac{\log L}{\log t}}$ they are deterministic functions of (Y, Y') . Thus when we use $W^{\frac{\log L}{\log t}}$ to extract $V^{\frac{\log L}{\log t}}$ from $X_{3^{\frac{\log L}{\log t}}}$, by Theorem 2.7 we have that

$$(\hat{V}, \hat{V}', Y, Y') \approx (U_m, \hat{V}', Y, Y'),$$

where the error is $O(\sum_{j=1}^{\frac{\log L}{\log t}} t^j \epsilon) + O((\frac{\log L}{\log t})^2 \epsilon) = O(L\epsilon)$. □

6 The Constructions of Non-Malleable Extractors

In this section we construct our improved seeded non-malleable extractors and seedless non-malleable extractors. Both the constructions follow the general approach developed in recent works [CGL16, CL16, Coh16a, Li17], i.e., first obtaining an advice and then applying an appropriate correlation breaker with advice. First we need the following advice generator from [CGL16].

Theorem 6.1 ([CGL16]). *There exist a constant $c > 0$ such that for all $n > 0$ and any $\epsilon > 0$, there exists an explicit function $\text{AdvGen} : \{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^L$ with $L = c \log(n/\epsilon)$ satisfying the following: Let X be an (n, k) -source, and Y be an independent uniform seed on d bits. Let Y' be a random variable on d bits s.t $Y' \neq Y$, and (Y, Y') is independent of X . Then with probability at least $1 - \epsilon$, $\text{AdvGen}(X, Y) \neq \text{AdvGen}(X, Y')$. Moreover, there is a deterministic function g such that $\text{AdvGen}(X, Y)$ is computed as follows. Let Y_1 be a small slice of Y with length $O(\log(n/\epsilon))$, compute $Z = \text{Ext}(X, Y_1)$ where Ext is an optimal seeded extractor from Theorem 2.7 which outputs $O(\log(n/\epsilon))$ bits. Finally compute $Y_2 = g(Y, Z)$ which outputs $O(\log(1/\epsilon))$ bits and let $\text{AdvGen}(X, Y) = (Y_1, Y_2)$.*

For two independent sources we also have the following slightly different advice generator.

Theorem 6.2 ([CGL16]). *There exist constants $0 < \gamma < \beta < 1$ such that for all $n > 0$ and any $\epsilon \geq \epsilon'$ for some $\epsilon' = 2^{-\Omega(n)}$, there exists an explicit function $\text{AdvGen} : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}^L$ with $L = 2\beta n + O(\log(1/\epsilon))$ satisfying the following: Let X, Y be two independent $(n, (1 - \gamma)n)$ -sources, and (X', Y') be some tampered versions of (X, Y) , such that (X, X') is independent of (Y, Y') . Furthermore either $X \neq X'$ or $Y \neq Y'$. Then with probability at least $1 - \epsilon$, $\text{AdvGen}(X, Y) \neq \text{AdvGen}(X', Y')$. Moreover, there is a deterministic function g such that $\text{AdvGen}(X, Y)$ is computed as follows. Let X_1, Y_1 be two small slice of X, Y respectively, with length βn , compute $Z = \text{IP}(X, Y_1)$ where IP is the inner product two source extractor from Theorem 2.8 which outputs $\Omega(n)$ bits. Finally compute $X_2 = g(X, Z), Y_2 = g(Y, Z)$ which both output $O(\log(1/\epsilon))$ bits and let $\text{AdvGen}(X, Y) = (X_1, X_2, Y_1, Y_2)$.*

By using these advice generators, the general approach of constructing seeded non-malleable extractors and seedless non-malleable extractors can be summarized in the following two theorems.

Theorem 6.3. [CGL16, CL16, Coh16a, Li17] *There is a constant $c > 1$ such that for any $n, k, d \in \mathbb{N}$ and $\epsilon_1, \epsilon_2 > 0$, if there is a $((k - c \log(n/\epsilon_1), d - c \log(n/\epsilon_1), \epsilon_2)$ advice correlation breaker $\text{AdvCB} : \{0, 1\}^k \times \{0, 1\}^d \times \{0, 1\}^{c \log(n/\epsilon_1)} \rightarrow \{0, 1\}^m$, then there exists an $(O(k), \epsilon_1 + \epsilon_2)$ seeded non-malleable extractor $\text{nmExt} : \{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^m$. Furthermore if $m \geq c \log(d/\epsilon_1)$ then there exists an $(O(k), \epsilon_1 + \epsilon_2)$ seeded non-malleable extractor $\text{nmExt} : \{0, 1\}^n \times \{0, 1\}^{O(d)} \rightarrow \{0, 1\}^{\Omega(k)}$.*

Sketch. The seeded non-malleable extractor is constructed as follows. First use the seed and the source to obtain an advice as in Theorem 6.1 with error $\epsilon_1/3$, however when we compute $Z = \text{Ext}(X, Y_1)$ we in fact output $Z_1 = \text{Ext}(X, Y_1)$ with k bits and choose Z to be a small slice of Z_1 with length $O(\log(n/\epsilon))$. Then we can fix the random variables $(Y_1, Y'_1, Z, Z', Y_2, Y'_2)$. Note that conditioned on this fixing (X, X') is still independent of (Y, Y') , and (Z_1, Z'_1) is a deterministic function of (X, X') thus is independent of (Y, Y') . Furthermore with probability $1 - \epsilon_1/3$, Z_1 has min-entropy at least $k - O(\log(n/\epsilon_1))$ and Y has min-entropy at least $d - O(\log(n/\epsilon_1))$. We can now apply the correlation breaker to (Z_1, Y) and the advice to get the desired output, where the total error is at most $\epsilon_1/3 + \epsilon_1/3 + \epsilon_1/3 + \epsilon_2 = \epsilon_1 + \epsilon_2$. If the output m is large enough (i.e., $m \geq c \log(d/\epsilon_1)$), then we can use it to extract from Y and then extract again from Z_1 to increase the output length to $\Omega(k)$. ■

Theorem 6.4. [CGL16, CL16, Coh16a, Li17] *There are constants $c > 1$, $0 < \gamma < \beta < 1/100$ such that for any $n \in \mathbb{N}$ and $\epsilon_1, \epsilon_2 > 0$, if there is a $((1 - 2\beta)n - c \log(n/\epsilon_1), (1 - 2\beta)n - c \log(n/\epsilon_1), \epsilon_2)$ advice correlation breaker $\text{AdvCB} : \{0, 1\}^n \times \{0, 1\}^n \times \{0, 1\}^{2\beta n + c \log(1/\epsilon_1)} \rightarrow \{0, 1\}^m$, then there exists an $((1 - \gamma)n, (1 - \gamma)n, \epsilon_1 + \epsilon_2)$ non-malleable two source extractor $\text{nmExt} : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}^m$. Furthermore if $m \geq c \log(n/\epsilon_1)$ then there exists an $((1 - \gamma)n, (1 - \gamma)n, \epsilon_1 + \epsilon_2)$ non-malleable two source extractor $\text{nmExt} : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}^{\Omega(n)}$.*

Sketch. The non-malleable two-source extractor is constructed as follows. First use the two independent sources (X, Y) to obtain an advice as in Theorem 6.2 with error $\epsilon_1/3$, then we can fix the random variables $(X_1, X'_1, Y_1, Y'_1, X_2, X'_2, Y_2, Y'_2)$. Note that conditioned on this fixing (X, X') is still independent of (Y, Y') , furthermore with probability $1 - \epsilon_1/3$, both X and Y have min-entropy at least $(1 - \gamma)n - \beta n - c \log(1/\epsilon_1) \geq (1 - 2\beta)n - c \log(1/\epsilon_1)$. We can now apply the correlation breaker to (X, Y) and the advice to get the desired output, where the total error is at most $\epsilon_1/3 + \epsilon_1/3 + \epsilon_1/3 + \epsilon_2 = \epsilon_1 + \epsilon_2$. If the output m is large enough (i.e., $m \geq c \log(d/\epsilon_1)$), then we can use it to extract from Y and then extract again from X to increase the output length to $\Omega(n)$. ■

Combined with our new correlation breakers with advice, we have the following new constructions of non-malleable extractors.

Theorem 6.5. *There exists a constant $C > 1$ such that for any constant $a \in \mathbb{N}, a \geq 2$, any $n, k \in \mathbb{N}$ and any $0 < \epsilon < 1$ with $k \geq C(\log n + a \log(1/\epsilon))$, there is an explicit construction of a strong seeded (k, ϵ) non-malleable extractor $\{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^m$ with $d = O(\log n) + \log(1/\epsilon)2^{O(a(\log \log(1/\epsilon))^{\frac{1}{a}})}$ and $m = \Omega(k)$. Alternatively, we can also achieve entropy $k \geq C \log n + \log(1/\epsilon)2^{C \cdot a(\log \log(1/\epsilon))^{\frac{1}{a}}}$ and $d = O(\log n + a \log(1/\epsilon))$.*

Proof. The theorem is obtained by combining Theorem 6.3, Lemma 5.4 and Lemma 4.9. We choose an error ϵ' to be the error in Theorem 6.3, Lemma 5.4 and Lemma 4.9. Thus the total error is $O(L\epsilon')$ where $L = O(\log(n/\epsilon'))$. To ensure $O(L\epsilon') = \epsilon$ it suffices to take $\epsilon' = \frac{\epsilon}{c \log(n/\epsilon)}$ for some constant $c > 1$. We know $\ell = O(\log(n/\epsilon'))$. Therefore to apply Lemma 5.4 and Lemma 4.9, we need to find m, d', d_1, d_2 such that

$$d' \geq c(\log \max\{d_1, d_2\} + \log(1/\epsilon')), m \geq d', d_1 \geq 8a \cdot d' + 6m \text{ and } d_2 \geq 8a \cdot d' + c^{a \cdot \log^{\frac{1}{a}} L} \cdot d'.$$

Then we can take

$$k = O(d_1 + m + \ell + \log(1/\epsilon')) \text{ and } d = O(d_2 + m + \ell + \log(1/\epsilon')).$$

It can be seen that we can take $m = O(\log(n/\epsilon'))$, $d' = O(\log \log n + \log(1/\epsilon'))$, $d_1 = 8a \cdot d' + 6m = O(\log n + a \log(1/\epsilon'))$ and $d_2 = 2^{O(a(\log \log(n/\epsilon'))^{\frac{1}{a}})} \cdot d'$. We now consider two cases. First, $\log(1/\epsilon') > \frac{\log n}{c^{a(\log \log n)^{\frac{1}{a}}}}$ for some large constant c' . In this case we have that

$$\log(1/\epsilon') > \frac{\log n}{c'^{a(\log \log n)^{\frac{1}{a}}}} > \sqrt{\log n}$$

for any $a \geq 2$. Thus

$$\log \log(n/\epsilon') = \log(\log n + \log(1/\epsilon')) < \log(\log^2(1/\epsilon') + \log(1/\epsilon')) < 2 \log \log(1/\epsilon') + 1.$$

Also note that $d' = O(\log \log n + \log(1/\epsilon')) = O(\log(1/\epsilon'))$. Thus in this case we have $d_2 \leq O(\log(1/\epsilon'))2^{O(a(\log \log(1/\epsilon'))^{\frac{1}{a}})} = \log(1/\epsilon')2^{O(a(\log \log(1/\epsilon'))^{\frac{1}{a}})}$. Next, consider the case where $\log(1/\epsilon') \leq \frac{\log n}{c'^{a(\log \log n)^{\frac{1}{a}}}}$. In this case note that we have $\log(1/\epsilon') < \log n$ and thus $2^{O(a(\log \log(n/\epsilon'))^{\frac{1}{a}})} < 2^{O(a(\log \log(n))^{\frac{1}{a}})}$. Therefore when c' is large enough and $a \geq 2$ we have that

$$d_2 \leq 2^{O(a(\log \log(n))^{\frac{1}{a}})}(\log \log n + \log(1/\epsilon')) \leq \log n.$$

Therefore altogether we have that $d_2 \leq (\log n + \log(1/\epsilon'))2^{O(a(\log \log(1/\epsilon'))^{\frac{1}{a}})}$ and $d = O(d_2 + m + \ell + \log(1/\epsilon')) = O(\log n) + \log(1/\epsilon')2^{O(a(\log \log(1/\epsilon'))^{\frac{1}{a}})}$. Note that $\log(1/\epsilon') = \log(1/\epsilon) + \log(\log n + \log(1/\epsilon)) + O(1)$, a careful analysis similar as above shows that we also have that

$$d = O(\log n) + \log(1/\epsilon)2^{O(a(\log \log(1/\epsilon))^{\frac{1}{a}})}.$$

Note that the correlation breaker is completely symmetric to both sources, and the only difference is in generating the advice. Thus after advice generation which costs both sources $O(\log(n/\epsilon))$ entropy, we can switch the role of the seed and the source. Therefore we can also get the other setting of parameters where $k \geq C \log n + \log(1/\epsilon) 2^{C \cdot a(\log \log(1/\epsilon))^{\frac{1}{a}}}$ and $d = O(\log n + a \log(1/\epsilon))$. ■

By using this theorem, we can actually improve the entropy requirement of the non-malleable extractor. Specifically, we have the following theorem.

Theorem 6.6. *There exists a constant $C > 1$ such that for any constant $a \in \mathbb{N}, a \geq 2$, any $n, k \in \mathbb{N}$ and any $0 < \epsilon < 1$ with $k \geq C(\log \log n + a \log(1/\epsilon))$, there is an explicit construction of a strong seeded (k, ϵ) non-malleable extractor $\{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^m$ with $d = O(\log n) + \log(1/\epsilon) 2^{O(a(\log \log(1/\epsilon))^{\frac{1}{a}})}$ and $m = \Omega(k)$. Alternatively, we can also achieve entropy $k \geq C \log \log n + \log(1/\epsilon) 2^{C \cdot a(\log \log(1/\epsilon))^{\frac{1}{a}}}$ and $d = O(\log n + a \log(1/\epsilon))$.*

Proof. We start by taking a slice of the seed Y_1 with length $O(\log(n/\epsilon))$ to extract from the source, and output some $k' = 0.9k$ uniform bits with error $\epsilon/2$. Note that conditioned on the fixing of (Y_1, Y_1') where Y_1' is the tampered version, the two sources are still independent, and the seed now has average conditional entropy at least $d - O(\log(n/\epsilon))$. We now switch the role of the seed and the source, and use the output of the extractor from the source as the seed of a non-malleable extractor and apply Theorem 6.5 with error $\epsilon/2$, so that the final error is ϵ .

Note that now we know the original seed is different from its tampered version, so we only need to obtain advice from the original seed and thus the advice size is $O(\log(d/\epsilon))$. Now we only need

$$k \geq C(\log d + a \log(1/\epsilon))$$

and

$$d - O(\log(n/\epsilon)) \geq C \log k + \log(1/\epsilon) 2^{C \cdot a(\log \log(1/\epsilon))^{\frac{1}{a}}}.$$

Thus we can choose

$$k \geq C'(\log \log n + a \log(1/\epsilon))$$

for some slightly larger constant $C' > 1$, while the requirement of the seed is still

$$d = O(\log n) + \log(1/\epsilon) 2^{O(a(\log \log(1/\epsilon))^{\frac{1}{a}})}.$$

Similarly, we can switch the role of the seed and the source to get the other setting of parameters. ■

The next theorem improves the seed length, at the price of using a slightly larger entropy.

Theorem 6.7. *There exists a constant $C > 1$ such that for any $n, k \in \mathbb{N}$ and $0 < \epsilon < 1$ with $k \geq C(\log n + \log(1/\epsilon) \log \log \log(1/\epsilon))$, there is an explicit construction of a strong seeded (k, ϵ) non-malleable extractor $\{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^m$ with $d = O(\log n + \log(1/\epsilon)(\log \log(1/\epsilon))^2)$ and $m = \Omega(k)$.*

Proof. The theorem is obtained by combining Theorem 6.3, Lemma 5.4 and Lemma 4.11. Again, We choose an error ϵ' to be the error in Theorem 6.3, Lemma 5.4 and Lemma 4.11. Thus the total error is $O(L\epsilon')$ where $L = O(\log(n/\epsilon'))$. To ensure $O(L\epsilon') = \epsilon$ it suffices to take $\epsilon' = \frac{\epsilon}{c \log(n/\epsilon)}$ for some constant $c > 1$. We also know $\ell = O(\log(n/\epsilon'))$ in Lemma 5.4. Thus to apply Lemma 4.11, we need to find m, d_1, d_2 such that (for simplicity, we choose $a = 4$ in Lemma 4.11),

$$m \geq c \log(d_2/\epsilon'), d_1 \geq 20c \log \log L \log(d_2/\epsilon') + 6m \text{ and } d_2 \geq 20c \log^2 L \log(d_1/\epsilon').$$

Then we can take

$$k = O(d_1 + m + \ell + \log(1/\epsilon')) \text{ and } d = O(d_2 + m + \ell + \log(1/\epsilon')).$$

A careful but tedious calculation shows that we can choose $k \geq C(\log n + \log(1/\epsilon') \log \log \log(1/\epsilon'))$ for some large enough constant $C > 1$, and $d = O(\log n + \log(1/\epsilon')(\log \log(1/\epsilon'))^2)$. Note that we can choose $m = O(\log(n/\epsilon'))$ for a large enough constant in $O(\cdot)$, thus by Theorem 6.3 we can get an output length of $\Omega(k)$. Finally, note that $\log(n/\epsilon') = O(\log(n/\epsilon))$, thus the theorem follows. ■

Similar to what we have done above, we can also use this to get improved parameters. Specifically, we have

Theorem 6.8. *There exists a constant $C > 1$ such that for any $n, k \in \mathbb{N}$ and $0 < \epsilon < 1$ with $k \geq C(\log \log n + \log(1/\epsilon) \log \log \log(1/\epsilon))$, there is an explicit construction of a strong seeded (k, ϵ) non-malleable extractor $\{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^m$ with $d = O(\log n + \log(1/\epsilon)(\log \log(1/\epsilon))^2)$ and $m = \Omega(k)$. Alternatively, we can also achieve entropy $k \geq C(\log \log n + \log(1/\epsilon)(\log \log(1/\epsilon))^2)$ and seed length $d = O(\log n + \log(1/\epsilon) \log \log(1/\epsilon))$.*

For non-malleable two-source extractors we have the following theorem.

Theorem 6.9. *There exists a constant $0 < \gamma < 1$ and a non-malleable two-source extractor for $(n, (1 - \gamma)n)$ sources with error $2^{-\Omega(n \log \log n / \log n)}$ and output length $\Omega(n)$.*

Proof. The theorem is obtained by combining Theorem 6.4, Lemma 5.6 and Lemma 4.9. Again, we choose an error ϵ' to be the error in Theorem 6.3, Lemma 5.4 and Lemma 4.11. Thus the total error is $O(L\epsilon')$ where $L = O(n)$. To ensure $O(L\epsilon') = \epsilon$ it suffices to take $\epsilon' = \frac{\epsilon}{cn}$ for some constant c . We also know $\ell = 2\beta n + o(n)$ for some constant $\beta < 1/100$ in Lemma 5.6. We choose $a = 2$ in Lemma 4.9 and thus we obtain a correlation breaker with $m = O(\log(n/\epsilon'))$, $d_1 = O(\log(n/\epsilon'))$ and $d_2 = \log(n/\epsilon')2^{O(\sqrt{\log t})}$ where t is the parameter in Construction 5.5 with $t \leq L$. Note that this also satisfies that $d_1 \geq 4m$ and $m \geq c \log(d_2/\epsilon)$ as required by Lemma 5.6.

Now we need to ensure that

$$(1 - \beta)n \geq c \frac{\log L}{\log t} \log(n/\epsilon') + \max\left\{8 \frac{\log L}{\log t} d_1, 2t \cdot d' + 4d_2\right\} + 5\ell + 4 \log(1/\epsilon'),$$

where $d' = O(\log(n/\epsilon'))$. We choose $t = \frac{\log L}{\log \log L}$ and this gives us

$$(1 - 12\beta)n \geq C \frac{\log L}{\log \log L} \log(n/\epsilon'),$$

for some constant $C > 1$. Note that $\log(n/\epsilon') = O(\log(n/\epsilon))$ thus we can set $\epsilon = 2^{-\Omega(n \log \log n / \log n)}$ and satisfy the above inequality. ■

For applications in two-source extractors, we first need the following generalization of non-malleable extractors, which allows multiple tampering.

Definition 6.10 (Seeded t -Non-malleable extractor). A function $\text{snmExt} : \{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^m$ is a seeded t -non-malleable extractor for min-entropy k and error ϵ if the following holds : If X is a source on $\{0, 1\}^n$ with min-entropy k and $\mathcal{A}_1, \dots, \mathcal{A}_t : \{0, 1\}^d \rightarrow \{0, 1\}^d$ are t arbitrary tampering functions with no fixed points, then

$$|\text{snmExt}(X, U_d) \circ \{\text{snmExt}(X, \mathcal{A}_i(U_d)), i \in [t]\} \circ U_d - U_m \circ \{\text{snmExt}(X, \mathcal{A}_i(U_d)), i \in [t]\} \circ U_d| < \epsilon$$

where U_m is independent of U_d and X .

The following theorem is a special case of Theorem 8.6 proved in [Li17].

Theorem 6.11. *Suppose there is a function f , a constant $\gamma > 0$ and an explicit non-malleable two-source extractor for $(f(\epsilon), (1 - \gamma)f(\epsilon))$ sources with error ϵ and output length $\Omega(f(\epsilon))$. Then there is a constant $C > 0$ such that for any $0 < \epsilon < 1$ with $k \geq Ct^2(\log n + f(\epsilon))$, there is an explicit strong seeded t -non-malleable extractor for (n, k) sources with seed length $d = Ct^2(\log n + f(\epsilon))$, error $O(t\epsilon)$ and output length $\Omega(f(\epsilon))$.*

Combined with Theorem 6.9, this immediately gives the following theorem.

Theorem 6.12. *There is a constant $C > 0$ such that for any $0 < \epsilon < 1$ and $n, k \in \mathbb{N}$ with $k \geq Ct^2(\log n + \frac{\log(1/\epsilon) \log \log(1/\epsilon)}{\log \log \log(1/\epsilon)})$, there is an explicit strong seeded t -non-malleable extractor for (n, k) sources with seed length $d = Ct^2(\log n + \frac{\log(1/\epsilon) \log \log(1/\epsilon)}{\log \log \log(1/\epsilon)})$, error $O(t\epsilon)$ and output length $\Omega(k/t^2)$. As a special case, there exists a seeded non-malleable extractor for entropy $k \geq C(\log n + \frac{\log(1/\epsilon) \log \log(1/\epsilon)}{\log \log \log(1/\epsilon)})$ and seed length $d = C(\log n + \frac{\log(1/\epsilon) \log \log(1/\epsilon)}{\log \log \log(1/\epsilon)})$.*

Similar techniques as above can reduce the $\log n$ term in the entropy requirement to $\log \log n$, so we get

Theorem 6.13. *There is a constant $C > 0$ such that for any $0 < \epsilon < 1$ and $n, k \in \mathbb{N}$ with $k \geq C(\log \log n + \frac{\log(1/\epsilon) \log \log(1/\epsilon)}{\log \log \log(1/\epsilon)})$, there is an explicit strong seeded non-malleable extractor for (n, k) sources with seed length and seed length $d = C(\log n + \frac{\log(1/\epsilon) \log \log(1/\epsilon)}{\log \log \log(1/\epsilon)})$.*

Ben-Aroya et. al [BADTS17] proved the following theorem.

Theorem 6.14. [BADTS17] *Suppose there is a function f and an explicit strong seeded t -non-malleable extractor (n, k') sources with seed length and entropy requirement $d = k' = f(t, \epsilon)$, then for every constant $\epsilon > 0$ there exist constants $t = t(\epsilon), c = c(\epsilon)$ and an explicit extractor $\text{Ext} : (\{0, 1\}^n)^2 \rightarrow \{0, 1\}$ for two independent (n, k) sources with $k \geq f(t, 1/n^c)$ and error ϵ .*

Combined with Theorem 6.7, this immediately gives the following theorem.

Theorem 6.15. *For every constant $\epsilon > 0$, there exists a constant $C > 1$ and an explicit two source extractor $\text{Ext} : (\{0, 1\}^n)^2 \rightarrow \{0, 1\}$ for entropy $k \geq C \frac{\log n \log \log n}{\log \log \log n}$ with error ϵ .*

7 Non-Malleable Two-Source Extractor and Non-Malleable Code

Formally, non-malleable codes are defined as follows.

Definition 7.1. [ADKO15] Let NM_k denote the set of trivial manipulation functions on k -bit strings, which consists of the identity function $I(x) = x$ and all constant functions $f_c(x) = c$, where $c \in \{0, 1\}^k$. Let $E : \{0, 1\}^k \rightarrow \{0, 1\}^m$ be an efficient randomized *encoding* function, and $D : \{0, 1\}^m \rightarrow \{0, 1\}^k$ be an efficient deterministic *decoding* function. Let $\mathcal{F} : \{0, 1\}^m \rightarrow \{0, 1\}^m$ be some class of functions. We say that the pair (E, D) defines an $(\mathcal{F}, k, \epsilon)$ -*non-malleable code*, if for all $f \in \mathcal{F}$ there exists a probability distribution G over NM_k , such that for all $x \in \{0, 1\}^k$, we have

$$|D(f(E(x))) - G(x)| \leq \epsilon.$$

Remark 7.2. The above definition is slightly different from the original definition in [DPW10]. However, [ADKO15] shows that the two definitions are equivalent.

We will mainly be focusing on the following family of tampering functions in this paper.

Definition 7.3. Given any $t > 1$, let \mathcal{S}_n^t denote the tampering family in the t -split-state-model, where the adversary applies t arbitrarily correlated functions h_1, \dots, h_t to t separate, n -bit parts of string. Each h_i can only be applied to the i -th part individually.

We remark that even though the functions h_1, \dots, h_t can be correlated, their correlation is independent of the original codewords. Thus, they are actually a convex combination of independent functions, applied to each part of the codeword. Therefore, without loss of generality we can assume that each h_i is a deterministic function, which acts on the i -th part of the codeword individually. We will mainly consider the case of $t = 2$, i.e., the two-split-state model. We recall the original definition of non-malleable two-source extractors by Cheraghchi and Guruswami [CG14b]. First we define the following function.

$$\text{copy}(x, y) = \begin{cases} x & \text{if } x \neq \text{same}^* \\ y & \text{if } x = \text{same}^* \end{cases}$$

Definition 7.4 (Seedless Non-Malleable 2-Source Extractor). A function $\text{nmExt} : (\{0, 1\}^n)^2 \rightarrow \{0, 1\}^m$ is a (k, ϵ) -seedless non-malleable extractor for two independent sources, if it satisfies the following property: Let X, Y be two independent (n, k) sources, and $f_1, f_2 : \{0, 1\}^n \rightarrow \{0, 1\}^n$ be two arbitrary tampering functions, then

1. $|\text{nmExt}(X, Y) - U_m| \leq \epsilon$.
2. There is a distribution \mathcal{D} over $\{0, 1\}^m \cup \{\text{same}^*\}$ such that for an independent Z sampled from \mathcal{D} , we have

$$(\text{nmExt}(X, Y), \text{nmExt}(f_1(X), f_2(Y))) \approx_\epsilon (\text{nmExt}(X, Y), \text{copy}(Z, \text{nmExt}(X, Y))).$$

Cheraghchi and Guruswami [CG14b] showed that the relaxed definition 1.5 implies the above general definition with a small loss in parameters. Specifically, we have

Lemma 7.5 ([CG14b]). *Let nmExt be a $(k - \log(1/\epsilon), \epsilon)$ -non-malleable two-source extractor according to Definition 1.5. Then nmExt is a $(k, 4\epsilon)$ -non-malleable two-source extractor according to Definition 7.4.*

The following theorem was proved by Cheraghchi and Guruswami [CG14b], which establishes a connection between seedless non-malleable extractors and non-malleable codes.

Theorem 7.6. *Let $\text{nmExt} : \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\}^m$ be a polynomial time computable seedless 2-non-malleable extractor at min-entropy n with error ϵ . Then there exists an explicit non-malleable code with an efficient decoder in the 2-split-state model with block length $= 2n$, rate $= \frac{m}{2n}$ and error $= 2^{m+1}\epsilon$.*

One can construct a non-malleable code in the 2-split-state model from a non-malleable two-source extractor as follows: Given any message $s \in \{0, 1\}^m$, the encoding $\text{Enc}(s)$ is done by outputting a uniformly random string from the set $\text{nmExt}^{-1}(s) \subset \{0, 1\}^{2n}$. Given any codeword $c \in \{0, 1\}^{2n}$, the decoding $\text{Dec}(c)$ is done by outputting $\text{nmExt}(c)$. Thus, to get an efficient encoder we need a way to efficiently uniformly sample from the pre-image of any output of the extractor.

Since our new non-malleable two-source extractor follows the same structure as in [Li17], we can use the same sampling procedure there to efficiently uniformly sample from the pre-image of any output of the extractor. We briefly recall the construction and sampling procedure in [Li17].

The extractor construction and sampling. The high level structure of the non-malleable two-source extractor in [Li17] is as follows. First take two small slices (X_1, Y_1) of both sources and apply the inner product based two-source extractor, as in Theorem 2.8. Then, use the output to sample $O(\log(1/\epsilon))$ bits from the encodings of both sources, using a randomness efficient sampler and an asymptotically good linear encoding of the sources. We need an asymptotically good encoding since then we only need to sample $O(\log(1/\epsilon))$ bits to ensure that the sampling of two different codewords are different with probability at least $1 - \epsilon$. The advice is then obtained by combining the slices and the sample bits. Now, take two larger slices (X_2, Y_2) of both sources and apply the correlation breaker. Finally, take another larger slice of either source (say X_3 from X) and apply a strong linear seeded extractor, which is easy to invert and has the same pre-image size for any output. By limiting the size of each slice to be small, the construction ensures that there are at least $n/2$ bits of each source that are only used in the encoding of the sources but never used in the subsequent extraction.

Now to sample uniformly from the pre-image of any output, we first uniformly independently generate the slices (X_1, Y_1, X_2, Y_2) and the sampled bits Z . From these we can compute the coordinates of the sampled bits and the output of the correlation breaker. Now we can invert the linear seeded extractor and uniformly sample X_3 given the output of the extractor and the output of the correlation breaker (which is used as the seed of the linear seeded extractor). Now, to sample the rest of the bits, we need to condition on the event that the sample bits from the encoding of the sources are indeed Z . Note that Z has size at most αn for some small constant $\alpha < 1/2$ since we can restrict the error to be at least some $2^{-\Omega(n)}$. Also note that for each source we have already sampled some bits but there are still at least $n/2$ un-sampled free bits, thus we insist on that no matter which αn columns of the generating matrix of the encoding we look at, the sub matrix corresponding to these columns and the last $n/2$ rows have full column rank. If this is true then no matter which coordinates we use and what Z is, the pre-image always have the same size and we can uniformly sample from the pre-image by solving a system of linear equations.

In [Li17], we use the Reed-Solomon encoding for each source with field \mathbb{F}_q for $q \approx n$. This is asymptotically good and also satisfies the property that any sub matrix with less columns than rows has full column rank since it is a Vandermonde matrix. However in this case each symbol has roughly $\log n$ bits so we can sample at most $n/\log n$ symbols (otherwise fixing them may already cost us all the entropy), thus the best error we can get using this encoding is $2^{-n/\log n}$. Here we need to get better error, so we use a binary generating matrix. It is easy to show using standard probabilist argument that there exists a binary generating matrix that satisfies our requirements.

Theorem 7.7. *There exists constants $0 < \alpha, \beta < 1$ such that for any $n \in \mathbb{N}$ there exists an $n \times m$ matrix over \mathbb{F}_2 with $n = \beta m$ which is the generating matrix of an asymptotically good code. Furthermore, Any sub-matrix formed by taking αn columns and the last $n/2$ rows has full column rank. In addition, for some $\epsilon = 2^{-O(n)}$, an ϵ -biased sample space over nm bits generates such a matrix with probability $1 - 2^{-\Omega(n)}$.*

Proof. We take an ϵ -biased sample space over nm bits for some $\epsilon = 2^{-O(n)}$. First, consider the sum of the rows over any non-empty subset of the rows. The sum is an m -bit string such that any non-empty parity is ϵ -close to uniform. Thus by the XOR lemma it is $2^{m/2}\epsilon$ -close to uniform. We know a uniform m -bit string has weight $d = m/4$ with probability at least $1 - 2^{-\Omega(m)}$. Thus for this string the probability is at least $1 - 2^{-\Omega(m)} - 2^{m/2}\epsilon$. By a union bound the total failure probability is at most $2^n(2^{-\Omega(m)} + 2^{m/2}\epsilon) = 2^{-\Omega(n)}$ by an appropriate choice of β and $\epsilon = 2^{-O(n)}$.

Next, consider any sub-matrix formed by taking βm columns and the last $n/2$ rows, if it's truly uniform, then the probability that it has full column rank is at least $1 - \alpha n 2^{\alpha n - n/2} \geq 1 - 2^{-n/4}$ for $\alpha < 1/5$. Now by a union bound the total failure probability is at most

$$\binom{m}{\alpha n} (2^{-n/4} + \epsilon) \leq \left(\frac{em}{\alpha n}\right)^{\alpha n} 2^{-n/4+1} = \left(\frac{e}{\beta\alpha}\right)^{\alpha n} 2^{-n/4+1},$$

if we choose $\epsilon < 2^{-n/4}$. Note that for a fixed β , the quantity $(\frac{e}{\beta\alpha})^\alpha$ goes to 1 as α goes to 0. Thus we can choose α small enough such that this failure probability is also $2^{-\Omega(n)}$. Therefore altogether the failure probability is $2^{-\Omega(n)}$. ■

Note that an ϵ -biased sample space over nm bits can be generated using $O(\log(nm/\epsilon)) = O(n)$ bits if $\epsilon = 2^{-O(n)}$. Now for any length $n \in \mathbb{N}$, we can compute the generating matrix (either using an ϵ -biased sample space or compute it deterministically in $2^{O(n)}$ time) once in the pre-processing step, and when we do encoding and decoding of the non-malleable code, all computation can be done in polynomial time.

Combining Theorem 7.6 and Theorem 6.9, we immediately obtain the following theorem.

Theorem 7.8. *For any $n \in \mathbb{N}$ there exists a non-malleable code with efficient encoder/decoder in the 2-split-state model with block length $2n$, rate $\Omega(\log \log n / \log n)$ and error $= 2^{-\Omega(n \log \log n / \log n)}$.*

8 Discussion and Open Problems

Several natural open problems remain here. The most intriguing one is how far we can push our new techniques. As mentioned above, one bottleneck here is that the computation of the merger is not a small space computation. If one can find a more succinct way to represent the computation, then it will certainly lead to further improvements (e.g., decrease the entropy requirement in two-source extractors to $O(\log n \sqrt{\log \log n})$). If in addition we can find a way to apply the recursive

construction as in Nisan’s generator [Nis92], then it is potentially possible to decrease the entropy requirement in two-source extractors to $O(\log n \log \log \log n)$. Finally, we believe our approach has the potential to eventually achieve truly optimal (up to constants) constructions.

References

- [ADKO15] D. Aggarwal, Y. Dodis, T. Kazana, and M. Obremski. Non-malleable reductions and applications. In *Proceedings of the 47th Annual ACM Symposium on Theory of Computing*, 2015.
- [ADL14] Divesh Aggarwal, Yevgeniy Dodis, and Shachar Lovett. Non-malleable codes from additive combinatorics. In *Proceedings of the 46th Annual ACM Symposium on Theory of Computing*, 2014.
- [Agg14] Divesh Aggarwal. Affine-evasive sets modulo a prime. Technical Report 2014/328, Cryptology ePrint Archive, 2014.
- [BADTS17] Avraham Ben-Aroya, Dean Doron, and Amnon Ta-Shma. Explicit two-source extractors for near-logarithmic min-entropy. In *Proceedings of the 49th Annual ACM Symposium on Theory of Computing*, 2017.
- [BBR88] Charles H. Bennett, Gilles Brassard, and Jean-Marc Robert. Privacy amplification by public discussion. *SIAM Journal on Computing*, 17(2):210–229, April 1988.
- [BIW04] Boaz Barak, R. Impagliazzo, and Avi Wigderson. Extracting randomness using few independent sources. In *Proceedings of the 45th Annual IEEE Symposium on Foundations of Computer Science*, pages 384–393, 2004.
- [BKS⁺05] Boaz Barak, Guy Kindler, Ronen Shaltiel, Benny Sudakov, and Avi Wigderson. Simulating independence: New constructions of condensers, Ramsey graphs, dispersers, and extractors. In *Proceedings of the 37th Annual ACM Symposium on Theory of Computing*, pages 1–10, 2005.
- [Bou05] Jean Bourgain. More on the sum-product phenomenon in prime fields and its applications. *International Journal of Number Theory*, 1:1–32, 2005.
- [BRSW06] Boaz Barak, Anup Rao, Ronen Shaltiel, and Avi Wigderson. 2 source dispersers for $n^{o(1)}$ entropy and Ramsey graphs beating the Frankl-Wilson construction. In *Proceedings of the 38th Annual ACM Symposium on Theory of Computing*, 2006.
- [CG88] Benny Chor and Oded Goldreich. Unbiased bits from sources of weak randomness and probabilistic communication complexity. *SIAM Journal on Computing*, 17(2):230–261, 1988.
- [CG14a] Mahdi Cheraghchi and Venkatesan Guruswami. Capacity of non-malleable codes. In *ITCS*, pages 155–168, 2014.
- [CG14b] Mahdi Cheraghchi and Venkatesan Guruswami. Non-malleable coding against bit-wise and split-state tampering. In *TCC*, pages 440–464, 2014.

- [CGL16] Eshan Chattopadhyay, Vipul Goyal, and Xin Li. Non-malleable extractors and codes, with their many tampered extensions. In *Proceedings of the 48th Annual ACM Symposium on Theory of Computing*, 2016.
- [CKOR10] N. Chandran, B. Kanukurthi, R. Ostrovsky, and L. Reyzin. Privacy amplification with asymptotically optimal entropy loss. In *Proceedings of the 42nd Annual ACM Symposium on Theory of Computing*, pages 785–794, 2010.
- [CL16] Eshan Chattopadhyay and Xin Li. Explicit non-malleable extractors, multi-source extractors and almost optimal privacy amplification protocols. In *Proceedings of the 57th Annual IEEE Symposium on Foundations of Computer Science*, 2016.
- [Coh15] Gil Cohen. Local correlation breakers and applications to three-source extractors and mergers. In *Proceedings of the 56th Annual IEEE Symposium on Foundations of Computer Science*, 2015.
- [Coh16a] Gil Cohen. Making the most of advice: New correlation breakers and their applications. In *Proceedings of the 57th Annual IEEE Symposium on Foundations of Computer Science*, 2016.
- [Coh16b] Gil Cohen. Non-malleable extractors - new tools and improved constructions. In *Proceedings of the 31st Annual IEEE Conference on Computational Complexity*, 2016.
- [Coh16c] Gil Cohen. Non-malleable extractors with logarithmic seeds. Technical Report TR16-030, ECCC, 2016.
- [Coh17] Gil Cohen. Two-source extractors for quasi-logarithmic min-entropy and improved privacy amplification protocols. In *Proceedings of the 49th Annual ACM Symposium on Theory of Computing*, 2017.
- [CRS14] Gil Cohen, Ran Raz, and Gil Segev. Non-malleable extractors with short seeds and applications to privacy amplification. *SIAM Journal on Computing*, 43(2):450–476, 2014.
- [CS16] Gil Cohen and Leonard Schulman. Extractors for near logarithmic min-entropy. In *Proceedings of the 57th Annual IEEE Symposium on Foundations of Computer Science*, 2016.
- [CZ14] Eshan Chattopadhyay and David Zuckerman. Non-malleable codes against constant split-state tampering. In *Proceedings of the 55th Annual IEEE Symposium on Foundations of Computer Science*, pages 306–315, 2014.
- [CZ16] Eshan Chattopadhyay and David Zuckerman. Explicit two-source extractors and resilient functions. In *Proceedings of the 48th Annual ACM Symposium on Theory of Computing*, 2016.
- [DKO13] Stefan Dziembowski, Tomasz Kazana, and Maciej Obremski. Non-malleable codes from two-source extractors. In *CRYPTO (2)*, pages 239–257, 2013.

- [DKRS06] Y. Dodis, J. Katz, L. Reyzin, and A. Smith. Robust fuzzy extractors and authenticated key agreement from close secrets. In *Advances in Cryptology — CRYPTO '06, 26th Annual International Cryptology Conference, Proceedings*, pages 232–250, 2006.
- [DKSS09] Zeev Dvir, Swastik Kopparty, Shubhangi Saraf, and Madhu Sudan. Extensions to the method of multiplicities, with applications to kakeya sets and mergers. In *Proceedings of the 50th Annual IEEE Symposium on Foundations of Computer Science*, 2009.
- [DLWZ14] Yevgeniy Dodis, Xin Li, Trevor D. Wooley, and David Zuckerman. Privacy amplification and non-malleable extractors via character sums. *SIAM Journal on Computing*, 43(2):800–830, 2014.
- [DORS08] Y. Dodis, R. Ostrovsky, L. Reyzin, and A. Smith. Fuzzy extractors: How to generate strong keys from biometrics and other noisy data. *SIAM Journal on Computing*, 38:97–139, 2008.
- [DP07] Stefan Dziembowski and Krzysztof Pietrzak. Intrusion-resilient secret sharing. In *Proceedings of the 48th Annual IEEE Symposium on Foundations of Computer Science, FOCS '07*, pages 227–237, Washington, DC, USA, 2007. IEEE Computer Society.
- [DPW10] Stefan Dziembowski, Krzysztof Pietrzak, and Daniel Wichs. Non-malleable codes. In *ICS*, pages 434–452, 2010.
- [DW08] Zeev Dvir and Avi Wigderson. Kakeya sets, new mergers and old extractors. In *Proceedings of the 49th Annual IEEE Symposium on Foundations of Computer Science*, 2008.
- [DW09] Yevgeniy Dodis and Daniel Wichs. Non-malleable extractors and symmetric key cryptography from weak secrets. In *Proceedings of the 41st Annual ACM Symposium on Theory of Computing*, pages 601–610, 2009.
- [Erd47] P. Erdős. Some remarks on the theory of graphs. *Bulletin of the American Mathematics Society*, 53:292–294, 1947.
- [GS18] Tom Gur and Igor Shinkar. An entropy lower bound for non-malleable extractors. Technical Report TR18-008, ECCO, 2018.
- [GUV09] Venkatesan Guruswami, Christopher Umans, and Salil Vadhan. Unbalanced expanders and randomness extractors from Parvaresh-Vardy codes. *Journal of the ACM*, 56(4), 2009.
- [KLR09] Yael Kalai, Xin Li, and Anup Rao. 2-source extractors under computational assumptions and cryptography with defective randomness. In *Proceedings of the 50th Annual IEEE Symposium on Foundations of Computer Science*, pages 617–628, 2009.
- [KLRZ08] Yael Tauman Kalai, Xin Li, Anup Rao, and David Zuckerman. Network extractor protocols. In *Proceedings of the 49th Annual IEEE Symposium on Foundations of Computer Science*, pages 654–663, 2008.

- [KOS17] Bhavana Kanukurthi, Lakshmi Bhavana Obbattu, and Sruthi Sekar. Four-state non-malleable codes with explicit constant rate. In *Fifteenth IACR Theory of Cryptography Conference*, 2017.
- [KR09] B. Kanukurthi and L. Reyzin. Key agreement from close secrets over unsecured channels. In *EUROCRYPT 2009, 28th Annual International Conference on the Theory and Applications of Cryptographic Techniques*, 2009.
- [Li11] Xin Li. Improved constructions of three source extractors. In *Proceedings of the 26th Annual IEEE Conference on Computational Complexity*, pages 126–136, 2011.
- [Li12a] Xin Li. Design extractors, non-malleable condensers and privacy amplification. In *Proceedings of the 44th Annual ACM Symposium on Theory of Computing*, pages 837–854, 2012.
- [Li12b] Xin Li. Non-malleable extractors, two-source extractors and privacy amplification. In *Proceedings of the 53rd Annual IEEE Symposium on Foundations of Computer Science*, pages 688–697, 2012.
- [Li13a] Xin Li. Extractors for a constant number of independent sources with polylogarithmic min-entropy. In *Proceedings of the 54th Annual IEEE Symposium on Foundations of Computer Science*, pages 100–109, 2013.
- [Li13b] Xin Li. New independent source extractors with exponential improvement. In *Proceedings of the 45th Annual ACM Symposium on Theory of Computing*, pages 783–792, 2013.
- [Li15a] Xin Li. Non-malleable condensers for arbitrary min-entropy, and almost optimal protocols for privacy amplification. In *12th IACR Theory of Cryptography Conference*, pages 502–531. Springer-Verlag, 2015. LNCS 9014.
- [Li15b] Xin Li. Three source extractors for polylogarithmic min-entropy. In *Proceedings of the 56th Annual IEEE Symposium on Foundations of Computer Science*, 2015.
- [Li16] Xin Li. Improved two-source extractors, and affine extractors for polylogarithmic entropy. In *Proceedings of the 57th Annual IEEE Symposium on Foundations of Computer Science*, 2016.
- [Li17] Xin Li. Improved non-malleable extractors, non-malleable codes and independent source extractors. In *Proceedings of the 49th Annual ACM Symposium on Theory of Computing*, 2017.
- [LRVW03] C. J. Lu, Omer Reingold, Salil Vadhan, and Avi Wigderson. Extractors: Optimal up to constant factors. In *Proceedings of the 35th Annual ACM Symposium on Theory of Computing*, pages 602–611, 2003.
- [Mek15] Raghu Meka. Explicit resilient functions matching Ajtai-Linial. *CoRR*, abs/1509.00092, 2015.

- [MW97] Ueli M. Maurer and Stefan Wolf. Privacy amplification secure against active adversaries. In *Advances in Cryptology — CRYPTO '97, 17th Annual International Cryptology Conference, Proceedings*, 1997.
- [Nis92] Noam Nisan. Pseudorandom generators for space-bounded computation. *Combinatorica*, 12:449–461, 1992.
- [NZ96] Noam Nisan and David Zuckerman. Randomness is linear in space. *Journal of Computer and System Sciences*, 52(1):43–52, 1996.
- [Rao06] Anup Rao. Extractors for a constant number of polynomially small min-entropy independent sources. In *Proceedings of the 38th Annual ACM Symposium on Theory of Computing*, 2006.
- [Raz05] Ran Raz. Extractors with weak random seeds. In *Proceedings of the 37th Annual ACM Symposium on Theory of Computing*, pages 11–20, 2005.
- [RW03] Renato Renner and Stefan Wolf. Unconditional authenticity and privacy from an arbitrarily weak secret. In *Advances in Cryptology — CRYPTO '03, 23rd Annual International Cryptology Conference, Proceedings*, pages 78–95, 2003.