

# Extractors and Lower Bounds for Locally Samplable Sources

Anindya De\*      Thomas Watson†

December 21, 2011

## Abstract

We consider the problem of extracting randomness from sources that are efficiently samplable, in the sense that each output bit of the sampler only depends on some small number  $d$  of the random input bits. As our main result, we construct a deterministic extractor that, given any  $d$ -local source with min-entropy  $k$  on  $n$  bits, extracts  $\Omega(k^2/nd)$  bits that are  $2^{-n^{\Omega(1)}}$ -close to uniform, provided  $d \leq o(\log n)$  and  $k \geq n^{2/3+\gamma}$  (for arbitrarily small constants  $\gamma > 0$ ).

Using our result, we also improve a result of Viola (FOCS 2010), who proved a  $1/2 - O(1/\log n)$  statistical distance lower bound for  $o(\log n)$ -local samplers trying to sample input-output pairs of an explicit boolean function, assuming the samplers use at most  $n + n^{1-\delta}$  random bits for some constant  $\delta > 0$ . Using a different function, we simultaneously improve the lower bound to  $1/2 - 2^{-n^{\Omega(1)}}$  and eliminate the restriction on the number of random bits.

## 1 Introduction

Randomness extraction is the following general problem. Given a sample from an imperfect physical source of randomness, which is modeled as a probability distribution on bit strings of length  $n$ , we wish to apply an efficient deterministic algorithm to the sample to produce an output which is almost uniformly distributed (and thus is suitable for use by a randomized algorithm). Of course, to extract randomness from a source, the source needs to “contain” a certain amount of randomness in the first place. It is well established that the most suitable measure of the amount of randomness in a source is its *min-entropy* (a distribution is said to have at least  $k$  bits of min-entropy if each outcome occurs with probability at most  $2^{-k}$ ). However, even if the source is known to have at least  $n - 1$  bits of min-entropy, no algorithm can extract even a single bit that is guaranteed to be close to uniformly distributed (see, for example, [Vad, Sha11] for proofs of this folklore observation). To deal with this problem, researchers have constructed *seeded extractors* (introduced by [NZ96]), which have access to a short uniformly random seed that is statistically independent of the source and which acts as a catalyst for the extraction process (see [Sha02, Vad, Sha11] for introductions).

However, there is a sense in which seeded extractors are overkill: They are guaranteed to work for completely arbitrary sources that have high enough min-entropy. It is reasonable to assume the physical source of randomness has some limited structure, in which case deterministic (that is,

---

\*Computer Science Division, University of California, Berkeley. This material is based upon work supported by the National Science Foundation under Grant No. CCF-1017403 and under Grant No. CCF-1118083.

†Computer Science Division, University of California, Berkeley. This material is based upon work supported by the National Science Foundation Graduate Research Fellowship under Grant No. DGE-0946797 and by the National Science Foundation under Grant No. CCF-1017403.

seedless) extraction may become viable. There are several classes of sources for which researchers have constructed good deterministic extractors. One such class is independent sources, where the  $n$  bits are partitioned into blocks which are assumed to be statistically independent of each other [CG88, DEOR04, BIW06, Bou05, BKS<sup>+</sup>10, Raz05, Sha08, Rao09a, BRSW06, Rao08, RZ08, RY11, TKLR09, Li11a]. Other such classes include so-called bit-fixing sources [CFG<sup>+</sup>85, KZ07, GRS06, Rao09b], affine sources [GR08a, Bou07, Rao09b, DG10, Yeh11, Li11b], polynomial sources [DGW09, BSG11], and algebraic varieties [Dvi09].

Trevisan and Vadhan [TV00] considered deterministic extractors for the class of sources that are samplable by efficient algorithms given uniform random bits. One may initially be concerned that extracting randomness from such sources is somehow circular or vacuous: We are assuming uniform random bits are used to sample the source, and our goal then is to “undo” the sampling and get uniform random bits back. The point is that this class of sources is just a model for physical sources. This is motivated by the following postulate about the universe: A physical source of randomness is generated by an efficient process in nature, so it is reasonable to model the source as being sampled by an efficient algorithm.

Trevisan and Vadhan constructed extractors for the class of sources samplable by general *time-bounded* algorithms, but their constructions are conditional on (somewhat non-standard) complexity-theoretic conjectures. It is common in other areas of research, such as proving lower bounds and constructing pseudorandom generators, that proving unconditional limits on the power of time-bounded algorithms is beyond the reach of current techniques. Thus researchers consider more restricted types of algorithms, such as small-space algorithms and bounded-depth circuits, which are combinatorially simple enough for us to prove unconditional results. Hence it is natural to try to construct unconditional deterministic extractors for sources samplable by such restricted algorithms. Kamp et al. [KRVZ11] succeeded in doing so for small-space samplers with streaming/one-way access to the random input bits.

However, at the time this paper was written, it was an open problem to construct an unconditional deterministic extractor for sources samplable by polynomial-size constant-depth circuits with unbounded fan-in gates. A basic obstacle is that this requires that input-output pairs of the extractor cannot be sampled by such circuits, and it was not even known how to construct an explicit function with the latter property. For example, although the parity function is known not to have subexponential-size constant-depth circuits [Yao85, Hås86], input-output pairs can be sampled very efficiently: Just take uniformly random bits  $x_1, \dots, x_n$  and output  $x_1, x_1 \oplus x_2, x_2 \oplus x_3, \dots, x_{n-1} \oplus x_n, x_n$ . In independent and concurrent work, Viola [Vio11] has constructed unconditional deterministic extractors for sources samplable by polynomial-size constant-depth circuits with unbounded fan-in gates, which in particular yields an explicit function whose input-output pairs cannot be sampled by such circuits (see Section 1.3).

Our goal in this paper is to expand the frontier of unconditional deterministic randomness extraction for sources with low-complexity samplers. We succeed in constructing extractors for sources samplable by small-depth circuits with *bounded* fan-in gates (which corresponds to the class  $\text{NC}^0$  when the depth is constant). This is equivalent to requiring that each output bit of the sampler only depends on a small number of input bits. We call such sources *locally samplable*. Even constructing extractors for sources where each output bit depends on at most one input bit is nontrivial, as such sources are a natural generalization of bit-fixing sources.

As pointed out above, a necessary condition for a function to be an extractor for sources sampled by a class of algorithms is that input-output pairs of the function cannot be sampled by

such algorithms. Finding explicit functions with the latter property is tougher than finding explicit functions that are hard to compute, because if a function is easy to compute, then input-output pairs can be obtained by just sampling a random input and then computing the corresponding output. Viola [Vio10] initiated the study of finding explicit boolean functions whose input-output pairs are hard to sample for low-complexity samplers (specifically, local samplers). Another contribution of our paper is an application of our extractor result to obtain an improvement of Viola’s result.

## 1.1 Results

We first give the formal definitions of extractors and locally samplable sources.

A distribution on a finite set  $S$  is said to have *min-entropy* at least  $k$  if each element of  $S$  occurs with probability at most  $2^{-k}$ . The *statistical distance* between two distributions  $D_1$  and  $D_2$  on a finite set  $S$  is defined to be  $\|D_1 - D_2\| = \max_{T \subseteq S} |\Pr_{D_1}[T] - \Pr_{D_2}[T]|$ . If  $\|D_1 - D_2\| \leq \epsilon$  then we also say  $D_1$  and  $D_2$  are  $\epsilon$ -close. If  $f : S \rightarrow S'$  and  $D$  is a distribution on  $S$ , then we let  $f(D)$  denote the distribution on  $S'$  obtained by drawing a sample from  $D$  and applying  $f$  to it. When we mention a distribution multiple times in an expression, all instantiations refer to a single sample from the distribution; for example,  $(D, f(D))$  denotes the distribution obtained by sampling  $w \sim D$  and outputting the pair  $(w, f(w))$ . We use  $U_n$  to denote the uniform distribution on  $\{0, 1\}^n$ . If  $\mathcal{C}$  is a class of distributions on  $\{0, 1\}^n$ , then a function  $\text{Ext} : \{0, 1\}^n \rightarrow \{0, 1\}^m$  is called a  $(k, \epsilon)$ -*extractor for  $\mathcal{C}$*  if for every distribution  $D \in \mathcal{C}$  with min-entropy at least  $k$ ,  $\|\text{Ext}(D) - U_m\| \leq \epsilon$ . Informally, when we say an extractor is *explicit* we mean that a uniform polynomial-time deterministic algorithm with the desired behavior is exhibited.

We define a  $d$ -*local sampler* to be a function  $f : \{0, 1\}^r \rightarrow \{0, 1\}^n$  such that each output bit depends on at most  $d$  input bits. In other words, for every  $j \in \{1, \dots, n\}$  there exists a subset  $I_j \subseteq \{1, \dots, r\}$  with  $|I_j| \leq d$  and a function  $f_j : \{0, 1\}^{|I_j|} \rightarrow \{0, 1\}$  such that the  $j^{\text{th}}$  output bit of  $f$  is obtained by evaluating  $f_j$  on the input bits indexed by  $I_j$ . The output distribution of the sampler is  $f(U_r)$ . We say a distribution  $D$  on  $\{0, 1\}^n$  is a  $d$ -*local source* if there exists a  $d$ -local sampler (with any input length  $r$ ) whose output distribution is  $D$ .

We have three main theorems. Our first main theorem gives an extractor for locally samplable sources.

**Theorem 1.** *For every constant  $\gamma > 0$  there exists a constant  $\beta > 0$  such that there exists an explicit  $(k, \epsilon)$ -extractor for the class of  $d$ -local sources with output length  $m = k^2/8nd$  and error  $\epsilon = 2^{-n^\beta}$ , provided  $k \geq n^{2/3+\gamma}$  and  $d \leq \beta \log n$ .*

Our second main theorem gives an extractor for 1-local sources (which generalize bit-fixing sources), achieving better min-entropy requirement and better output length than Theorem 1.

**Theorem 2.** *For every constant  $\gamma > 0$  there exists a constant  $\beta > 0$  such that there exists an explicit  $(k, \epsilon)$ -extractor for the class of 1-local sources with output length  $m = k - o(k)$  and error  $\epsilon = 2^{-n^\beta}$ , provided  $k \geq n^{1/2+\gamma}$ .*

Our third main theorem concerns the problem of finding explicit functions whose input-output pairs are hard to sample, as discussed in the paragraph right before Section 1.1.

**Theorem 3.** *There exists a universal constant  $\beta > 0$  and an explicit function  $F : \{0, 1\}^n \rightarrow \{0, 1\}$  such that for every  $d$ -local source  $D$  on  $\{0, 1\}^{n+1}$  with  $d \leq \beta \log n$ ,  $\|D - (U_n, F(U_n))\| \geq 1/2 - 2^{-n^\beta}$ .*

## 1.2 Techniques

We now discuss the techniques we use to prove these three theorems. The proof of Theorem 1 has three steps.

The first step is to construct a certain extractor for 1-local sources (which in particular yields Theorem 2). To do this, we observe that extractors for so-called low-weight affine sources also work for 1-local sources. Rao [Rao09b] constructed an extractor for low-weight affine sources. Using Rao’s extractor off-the-shelf would lead to a weaker version of Theorem 1 with min-entropy requirement  $k \geq n^{1-\gamma}$  for *some* constant  $\gamma > 0$ . To improve the min-entropy requirement, we construct an improved extractor for low-weight affine sources by building on [Rao09b]. While Rao’s extractor handles affine sources of min-entropy at least  $k$  and weight at most  $k^\gamma$  for some constant  $\gamma > 0$ , our improvement handles sources with weight at most  $k^{1-\gamma}$  for *any* constant  $\gamma > 0$ . The key ingredient in our improvement is the strong condenser of Guruswami, Umans, and Vadhan [GUV09]. We present this step in Section 3 and Section 7.

The second step is to show that extractors for 1-local sources also work for  $o(\log n)$ -local sources. To do this, we relate the problem to a concept we call *superindependent matchings* in bipartite graphs, and we prove a combinatorial lemma about the existence of such matchings. We present this step in Section 4.

The third step is to increase the output length of the extractor using the technique of “obtaining an independent seed” introduced by Gabizon et al. [GRS06] (see also [Sha08]). Combining step 1 and step 2 yields an extractor with output length  $\Omega(k^2/nd^{32^d})$ . To increase the output length to  $\Omega(k^2/nd)$ , we adapt the technique from [GRS06]. A key ingredient in our argument is a lemma due to Vadhan [Vad04], which is a strengthened version of a classic lemma due to Nisan and Zuckerman [NZ96]. While the result of [GRS06] achieves output length  $k - o(k)$  for bit-fixing sources, we lose a factor of  $\Omega(k/n)$  in the output length due to the way we use Vadhan’s lemma, and we lose another factor of  $\Omega(1/d)$  since conditioning on  $p$  bits of the output of a  $d$ -local sampler could cause a loss of  $pd$  bits of min-entropy. We present this step in Section 5.

Viola [Vio10] proved a version of Theorem 3 where the statistical distance lower bound is only  $1/2 - O(1/\log n)$ , and the  $d$ -local sampler is restricted to use at most  $n + n^{1-\delta}$  random bits for any constant  $\delta > 0$ . His function  $F$  is what he calls “majority mod  $p$ ”. Using a different function  $F$  (namely, any bit of the extractor underlying Theorem 1), we simultaneously improve the lower bound to  $1/2 - 2^{-n^{\Omega(1)}}$  and eliminate the restriction on the number of random bits. Our proof of Theorem 3 uses ideas similar to Viola’s, but is actually somewhat simpler given the extraction property of  $F$ . In [Vio10], Viola also showed that for symmetric functions  $F$ , one cannot hope to get such a strong lower bound for samplers that are polynomial-size constant-depth circuits. Our extractor function  $F$  is not symmetric. We present the proof of Theorem 3 in Section 6.<sup>1</sup>

## 1.3 Concurrent Work

In independent and concurrent work, Viola [Vio11] obtained extractors for  $d$ -local sources with  $d \leq n^{o(1)}$  and for sources sampled by polynomial-size constant-depth circuits. The high level idea behind the extractor is the same as in our work: Show that the given source is close to a convex

---

<sup>1</sup>We also mention in passing that Lovett and Viola [LV11] exhibited an explicit distribution on  $\{0, 1\}^n$  that cannot be sampled within statistical distance  $1 - 1/n^{\Omega(1)}$  by polynomial-size constant-depth circuits, namely the uniform distribution over the codewords of any asymptotically good error-correcting code. However, this distribution is not of the same form as sampling input-output pairs.

combination of 1-local sources, and use the extractor in [Rao09b]. However, the proofs in [Vio11] are much more involved than in this paper. For  $d$ -local sources with  $d \leq n^{o(1)}$ , Viola requires min-entropy  $k \geq n^{3/4+\gamma}$  (for any constant  $\gamma > 0$ ) and achieves output length  $m = \tilde{\Omega}(k^3/n^2d^3)$  and error  $\epsilon = 2^{-n^{\Omega(1)}}$  (though the output length can be improved to  $\Omega(k^2/nd)$  using the technique we present in Section 5 based on [GRS06]). When  $d \leq o(\log n)$  he obtains a result similar to our Theorem 1 but with worse output length: He requires min-entropy  $k \geq n^{2/3+\gamma}$  and achieves output length  $m = \Omega(k^2/nd^22^d)$  and error  $\epsilon = 2^{-n^{\Omega(1)}}$ . For sources sampled by polynomial-size constant-depth circuits, he requires min-entropy  $k \geq n^{2/3+\gamma}$  and achieves output length  $m = \Omega(k^2/n^{1+\Omega(1)})$  and error  $\epsilon = n^{-\omega(1)}$ .

## 1.4 Previous Work on the Power of Locally Computable Functions

There has been a substantial amount of work on whether various cryptographic and complexity-theoretic objects can be computed locally. Several works [DM04, Lu04, Vad04, AIK08, Zim10, DT09] have studied the problem of constructing locally computable seeded extractors (that is, the extractor itself is locally computable, as opposed to our setting where the sampler for the source is locally computable). A variety of works [CM01, MST06, AIK06, AIK08, IKOS08, App11, ABR11] have given positive and negative results on the existence of locally computable pseudorandom generators. Several works [Hås87, Gol11, AIK06, CEMT09, BQ09] have studied the possibility of locally computable one-way functions. Goldwasser et al. [GGH<sup>+</sup>07] gave positive and negative results on interactive proof systems with locally computable verifiers. Arora et al. [ASW09] show that the adjacency list of certain logarithmic-degree expander graphs can be computed with constant locality, and they ask whether the same holds for constant-degree expander graphs.

## 2 Preliminaries

In this paper we work with bipartite graphs  $G = (L, R, E)$ , where  $L, R$  are disjoint finite sets (the left and right nodes) and  $E$  is a set of unordered pairs where one element comes from  $L$  and the other from  $R$ . The distance between two nodes is the number of edges on a shortest path between them.

To every function  $f : \{0, 1\}^r \rightarrow \{0, 1\}^n$  we associate a bipartite graph  $G = (L, R, E)$  where  $L = \{1, \dots, r\} \times \{\text{in}\}$ ,  $R = \{1, \dots, n\} \times \{\text{out}\}$ , and  $\{(i, \text{in}), (j, \text{out})\} \in E$  if and only if the  $j^{\text{th}}$  output bit of  $f$  depends on the  $i^{\text{th}}$  input bit of  $f$  (that is, for some setting of all input bits except the  $i^{\text{th}}$ , the  $j^{\text{th}}$  output bit equals the  $i^{\text{th}}$  input bit or its complement). Note that we include no unnecessary edges, and the graph is unique. We use  $I_j \times \{\text{in}\}$  to denote the set of neighbors of  $(j, \text{out})$  and  $J_i \times \{\text{out}\}$  to denote the set of neighbors of  $(i, \text{in})$ . Observe that if  $f(U_r)$  has min-entropy at least  $k$ , then there are at least  $k$  non-isolated nodes in  $L$ , and in particular  $r \geq k$ .

We say  $f$  is a  $d$ -local sampler if each node in  $R$  has degree at most  $d$ , and we say a distribution on  $\{0, 1\}^n$  is a  $d$ -local source if it equals  $f(U_r)$  for some  $d$ -local sampler  $f$  (with any input length  $r$ ). We say  $f$  is a  $(d, c)$ -local sampler if each node in  $R$  has degree at most  $d$  and each node in  $L$  has degree at most  $c$ , and we say a distribution on  $\{0, 1\}^n$  is a  $(d, c)$ -local source if it equals  $f(U_r)$  for some  $(d, c)$ -local sampler  $f$  (with any input length  $r$ ).

Suppose  $Y$  is a finite set of indices,  $(p_y)_{y \in Y}$  is a distribution on  $Y$ , and for each  $y \in Y$ ,  $D_y$  is a distribution on a finite set  $S$ . Then the *convex combination*  $\sum_{y \in Y} p_y D_y$  is defined to be the distribution on  $S$  obtained by sampling  $y$  according to  $(p_y)_{y \in Y}$ , then outputting a sample from  $D_y$ .

**Lemma 1.** Suppose  $\text{Ext} : \{0, 1\}^n \rightarrow \{0, 1\}^m$  is any function and  $D = \sum_{y \in Y} p_y D_y$  is a distribution on  $\{0, 1\}^n$ . Then for every  $\epsilon \geq 0$ ,

$$\|\text{Ext}(D) - U_m\| \leq \epsilon + \Pr_{y \sim (p_y)_{y \in Y}} \left[ \|\text{Ext}(D_y) - U_m\| > \epsilon \right].$$

*Proof.* First, observe that  $\text{Ext}(D) = \sum_{y \in Y} p_y \text{Ext}(D_y)$ . Now for every  $T \subseteq \{0, 1\}^n$  we have

$$\begin{aligned} & \left| \Pr_{\text{Ext}(D)}[T] - \Pr_{U_m}[T] \right| \\ &= \left| \sum_{y \in Y} p_y (\Pr_{\text{Ext}(D_y)}[T] - \Pr_{U_m}[T]) \right| \\ &\leq \sum_{y \in Y} p_y \left| \Pr_{\text{Ext}(D_y)}[T] - \Pr_{U_m}[T] \right| \\ &\leq \epsilon \cdot \Pr_{y \sim (p_y)_{y \in Y}} \left[ \left| \Pr_{\text{Ext}(D_y)}[T] - \Pr_{U_m}[T] \right| \leq \epsilon \right] + 1 \cdot \Pr_{y \sim (p_y)_{y \in Y}} \left[ \|\text{Ext}(D_y) - U_m\| > \epsilon \right] \\ &\leq \epsilon + \Pr_{y \sim (p_y)_{y \in Y}} \left[ \|\text{Ext}(D_y) - U_m\| > \epsilon \right] \end{aligned}$$

which gives the desired bound on  $\|\text{Ext}(D) - U_m\|$ .  $\square$

**Corollary 1.** Suppose every distribution in  $\mathcal{C}$  with min-entropy at least  $k$  can be written as a convex combination  $\sum_{y \in Y} p_y D_y$  where  $\Pr_{y \sim (p_y)_{y \in Y}} [D_y \text{ is in } \mathcal{C}' \text{ and has min-entropy at least } k'] \geq 1 - \delta$ . Then every  $(k', \epsilon')$ -extractor for  $\mathcal{C}'$  is also a  $(k, \epsilon)$ -extractor for  $\mathcal{C}$  where  $\epsilon = \epsilon' + \delta$ .

**Corollary 2.** Suppose every distribution in  $\mathcal{C}$  with min-entropy at least  $k$  is a convex combination of distributions in  $\mathcal{C}'$  with min-entropy at least  $k'$ . Then every  $(k', \epsilon)$ -extractor for  $\mathcal{C}'$  is also a  $(k, \epsilon)$ -extractor for  $\mathcal{C}$ .

**Lemma 2.** Every  $d$ -local source with min-entropy at least  $k$  is a convex combination of  $(d, c)$ -local sources with min-entropy at least  $k - nd/c$ .

*Proof.* Consider an arbitrary  $d$ -local sampler  $f : \{0, 1\}^r \rightarrow \{0, 1\}^n$  whose output distribution has min-entropy at least  $k$ , and let  $G = (L, R, E)$  be the associated bipartite graph. Since  $|E| \leq nd$ , there are at most  $nd/c$  nodes in  $L$  with degree greater than  $c$ ; without loss of generality these nodes are  $\{r - \ell + 1, \dots, r\} \times \{\text{in}\}$  for some  $\ell \leq nd/c$ . For each string  $y \in \{0, 1\}^\ell$ , define  $f_y : \{0, 1\}^{r-\ell} \rightarrow \{0, 1\}^n$  as  $f_y(x) = f(x, y)$  (hardwiring the last  $\ell$  bits to  $y$ ). Then  $f(U_r) = \sum_{y \in \{0, 1\}^\ell} \frac{1}{2^\ell} f_y(U_{r-\ell})$ . Moreover, each  $f_y(U_{r-\ell})$  is a  $(d, c)$ -local source with min-entropy at least  $k - nd/c$ , since if some  $z \in \{0, 1\}^n$  and  $y^* \in \{0, 1\}^\ell$  satisfied  $\Pr_{x \sim U_{r-\ell}} [f_{y^*}(x) = z] > 1/2^{k-nd/c}$  then we would have

$$\Pr_{x \sim U_{r-\ell}, y \sim U_\ell} [f(x, y) = z] \geq \Pr_{y \sim U_\ell} [y = y^*] \cdot \Pr_{x \sim U_{r-\ell}} [f(x, y^*) = z] > \frac{1}{2^\ell} \cdot \frac{1}{2^{k-nd/c}} \geq 1/2^k$$

contradicting that  $f(U_r)$  has min-entropy at least  $k$ .  $\square$

In this paper we also make use of seeded extractors. A function  $\text{Ext} : \{0, 1\}^n \times \{0, 1\}^t \rightarrow \{0, 1\}^m$  is called a *seeded*  $(k, \epsilon)$ -extractor if for every distribution  $D$  on  $\{0, 1\}^n$  with min-entropy at least  $k$ ,  $\|\text{Ext}(D, U_t) - U_m\| \leq \epsilon$  where  $U_t$  is independent of  $D$ . We say  $\text{Ext}$  is a *strong seeded*  $(k, \epsilon)$ -extractor if for every distribution  $D$  on  $\{0, 1\}^n$  with min-entropy at least  $k$ ,<sup>2</sup>

$$\Pr_{y \sim U_t} \left[ \|\text{Ext}(D, y) - U_m\| \leq \epsilon \right] \geq 1 - \epsilon.$$

<sup>2</sup>According to this definition, every strong seeded  $(k, \epsilon)$ -extractor is also a seeded  $(k, 2\epsilon)$ -extractor.

We say Ext is *linear* if for every seed  $y \in \{0, 1\}^t$ , the function  $\text{Ext}(\cdot, y) : \{0, 1\}^n \rightarrow \{0, 1\}^m$  is linear over  $\mathbb{F}_2$ , where  $\mathbb{F}_q$  denotes the finite field of size  $q$ .

If  $z \in \{0, 1\}^n$  and  $J \subseteq \{1, \dots, n\}$ , then we let  $z|_J \in \{0, 1\}^{|J|}$  denote the substring of  $z$  indexed by the coordinates in  $J$ . If  $D$  is a distribution on  $\{0, 1\}^n$  and  $J \subseteq \{1, \dots, n\}$ , then we let  $D|_J$  denote the marginal distribution on the coordinates in  $J$ .

Finally, all logarithms in this paper are base 2.

### 3 1-Local Sources

An *affine source* is a distribution on  $\{0, 1\}^n$  which is uniform over an affine subspace (where  $\{0, 1\}^n$  is viewed as a vector space over  $\mathbb{F}_2$ ). If the subspace has dimension  $k$  then it has size  $2^k$  and hence the source has min-entropy  $k$ . The distribution can be sampled by picking  $x_1, \dots, x_k \in \{0, 1\}$  uniformly at random and outputting  $z_0 + x_1 z_1 + \dots + x_k z_k$  where  $z_0 \in \{0, 1\}^n$  is a shift vector and  $z_1, \dots, z_k \in \{0, 1\}^n$  are a basis of the associated linear subspace. The source is said to be a *weight- $c$  affine source* if there exist basis vectors  $z_1, \dots, z_k$  each of which has Hamming weight at most  $c$ .

**Observation 1.** *Every  $(1, c)$ -local source is also a weight- $c$  affine source.*

*Proof.* Consider an arbitrary  $(1, c)$ -local sampler  $f : \{0, 1\}^k \rightarrow \{0, 1\}^n$  and assume without loss of generality that there are no isolated nodes on the left side of the associated bipartite graph. For each  $i \in \{1, \dots, k\}$ , let  $J_i \times \{\text{out}\}$  be the set of neighbors of  $(i, \text{in})$ , and let  $1_{J_i} \in \{0, 1\}^n$  be the characteristic vector of this set. For each  $i \in \{1, \dots, k\}$  we have  $|J_i| \leq c$  and hence  $1_{J_i}$  has Hamming weight at most  $c$  (since  $f$  is a  $(1, c)$ -local sampler). It is straightforward to verify that the output distribution of  $f$  is sampled by picking  $x_1, \dots, x_k \in \{0, 1\}$  uniformly at random and outputting  $f(0^k) + x_1 1_{J_1} + \dots + x_k 1_{J_k}$ . Moreover, the vectors  $1_{J_i}$  are linearly independent.  $\square$

Rao [Rao09b] (building on [Rao09a]) constructed extractors for low-weight affine sources.

**Theorem 4 ([Rao09b]).** *There exist universal constants  $C, \gamma > 0$  such that for all  $k \geq \log^C n$  there exists an explicit  $(k, 2^{-k^{\Omega(1)}})$ -extractor with output length  $m = k - o(k)$  for the class of weight- $k^\gamma$  affine (and in particular,  $(1, k^\gamma)$ -local) sources.*

We improve Rao's result to obtain the following theorem, which we prove in Section 7.

**Theorem 5.** *There exists a universal constant  $C > 0$  such that for every constant  $\gamma > 0$  and all  $k \geq \log^{C/\gamma} n$  there exists an explicit  $(k, 2^{-k^{\Omega(1)}})$ -extractor with output length  $m = k - o(k)$  for the class of weight- $k^{1-\gamma}$  affine (and in particular,  $(1, k^{1-\gamma})$ -local) sources.*

We now explain how Theorem 2 follows from Theorem 5, Lemma 2, and Corollary 2. We first note the following immediate corollary of Theorem 5.

**Corollary 3.** *For every constant  $\gamma > 0$  there exists a constant  $\beta > 0$  such that for all  $k \geq n^{1/2+\gamma}$  there exists an explicit  $(k, 2^{-n^\beta})$ -extractor with output length  $m = k - o(k)$  for the class of weight- $n^{1/2}$  affine (and in particular,  $(1, n^{1/2})$ -local) sources.*

Lemma 2 implies that every 1-local source with min-entropy at least  $k \geq n^{1/2+\gamma}$  is a convex combination of  $(1, n^{1/2})$ -local sources with min-entropy at least  $k - n^{1/2} \geq k - o(k)$ . Theorem 2 then follows from Corollary 2 and Corollary 3.

Bourgain [Bou07], Yehudayoff [Yeh11], and Li [Li11b] constructed extractors for linear min-entropy affine sources (of arbitrary weight), achieving better error but worse output length than Theorem 5.

**Theorem 6 ([Bou07]).** *For every constant  $\delta > 0$  there exists an explicit  $(\delta n, 2^{-\Omega(n)})$ -extractor with output length  $m = \Omega(n)$  for the class of affine (and in particular, 1-local) sources.*

Theorem 6 can be used to improve the error in Theorem 1 and Theorem 2 when  $k \geq \Omega(n)$  and  $d \leq O(1)$ . We omit the details, so as to avoid having a laundry list of results.

## 4 $d$ -Local Sources

The following theorem shows that to get extractors for  $d$ -local sources, it suffices to construct extractors for 1-local sources.

**Theorem 7.** *Every  $(k', \epsilon')$ -extractor for  $(1, 2nd/k)$ -local sources is also a  $(k, \epsilon)$ -extractor for  $d$ -local sources, where  $k' = k^2/4nd^32^d$  and  $\epsilon = \epsilon' + e^{-k'/4}$ .*

Assuming  $k \geq n^{2/3+\gamma}$  (for constant  $\gamma > 0$ ) and  $d \leq \beta \log n$  (for small enough constant  $\beta > 0$ ) in Theorem 7, we find that it suffices to have a  $(k', \epsilon')$ -extractor for  $(1, c)$ -local sources where  $k' \geq n^{1/3+\gamma}$  and  $c = 2nd/k \leq n^{1/3} \leq (k')^{1-\gamma}$ . Such an extractor is given by Theorem 5, with error  $\epsilon' = 2^{-n^{\Omega(1)}}$  (and thus  $\epsilon = \epsilon' + e^{-k'/4} \leq 2^{-n^{\Omega(1)}}$ ). This already yields a version of Theorem 1 with output length  $k' - o(k') = \Omega(k^2/nd^32^d)$ .

As a corollary to Theorem 7, we also find that if we could construct an explicit extractor for 1-local sources with min-entropy at least  $n^\gamma$  for arbitrarily small constants  $\gamma > 0$  (with output length  $m \geq 1$  and error  $\epsilon \leq 1/2$ , say) then we would get explicit extractors for  $o(\log n)$ -local sources with min-entropy at least  $n^{1/2+\gamma}$  for arbitrarily small constants  $\gamma > 0$ . This  $n^{1/2}$  min-entropy barrier is common in extractor constructions.

### 4.1 Superindependent Matchings

We first prove a combinatorial lemma that is needed for the proof of Theorem 7.

**Definition 1.** *Given a bipartite graph  $G = (L, R, E)$ , we say a set of edges  $M \subseteq E$  is a superindependent matching if there is no path of length at most two in  $G$  from an endpoint of an edge in  $M$  to an endpoint of a different edge in  $M$ .*

**Lemma 3.** *Suppose  $G = (L, R, E)$  is a bipartite graph with no isolated nodes and such that each node in  $L$  has degree at most  $c$  and each node in  $R$  has degree at most  $d$ . Then  $G$  has a superindependent matching of size at least  $|L|/d^2c$ .*

*Proof.* Let  $M$  be a largest superindependent matching in  $G$ , and suppose for contradiction that  $|M| < |L|/d^2c$ . Note that for each node in  $R$ , the number of nodes in  $L$  within distance three in  $G$  is at most  $d(1 + (c-1)(d-1)) \leq d^2c$ . Thus the number of nodes in  $L$  within distance three of the right endpoints of edges in  $M$  is at most  $|M| \cdot d^2c < |L|$ . Hence there exists a node  $u \in L$  at distance greater than three from the right endpoint of every edge in  $M$ . Since  $G$  has no isolated nodes, there exists a node  $v \in R$  such that  $\{u, v\} \in E$ . Note that there is no path of length at most two from either  $u$  or  $v$  to an endpoint of an edge in  $M$ , since otherwise a simple case analysis would show that  $u$  is within distance three of the right endpoint of an edge in  $M$ . Thus  $M \cup \{\{u, v\}\}$  is a superindependent matching, contradicting the maximality of  $M$ .  $\square$



## 4.2 Proof of Theorem 7

Suppose  $\text{Ext} : \{0, 1\}^n \rightarrow \{0, 1\}^m$  is a  $(k', \epsilon')$ -extractor for  $(1, 2nd/k)$ -local sources. By Corollary 2 and Lemma 2 it suffices to show that  $\text{Ext}$  is a  $(k/2, \epsilon)$ -extractor for  $(d, c)$ -local sources where  $c = 2nd/k$ . The plan is to show that every  $(d, c)$ -local source with min-entropy at least  $k/2$  is a convex combination of  $(1, c)$ -local sources most of which have min-entropy at least  $k'$ , and then apply Corollary 1.

So consider an arbitrary  $(d, c)$ -local sampler  $f : \{0, 1\}^r \rightarrow \{0, 1\}^n$  whose output distribution has min-entropy at least  $k/2$ , and let  $G = (L, R, E)$  be the associated bipartite graph. If we obtain  $\tilde{G}$  from  $G$  by removing any isolated nodes, then  $\tilde{G}$  still has at least  $k/2$  nodes on its left side. Applying Lemma 3 to  $\tilde{G}$  tells us that  $G$  has a superindependent matching  $M$  of size at least  $k/(2d^2c)$ . Let  $\ell = |M|$ , and without loss of generality assume that the left endpoints of  $M$  are  $L' = \{1, \dots, \ell\} \times \{\text{in}\}$ . We write inputs to  $f$  as  $(x, y)$  where  $x \in \{0, 1\}^\ell$  and  $y \in \{0, 1\}^{r-\ell}$ . Since  $M$  is superindependent, each node in  $R$  is adjacent to at most one node in  $L'$ . Thus if we define  $f_y : \{0, 1\}^\ell \rightarrow \{0, 1\}^n$  as  $f_y(x) = f(x, y)$  (hardwiring the last  $r - \ell$  input bits to  $y$ ) then for each  $y$ ,  $f_y$  is a  $(1, c)$ -local sampler. Observe that  $f(U_r) = \sum_{y \in \{0, 1\}^{r-\ell}} \frac{1}{2^{r-\ell}} f_y(U_\ell)$ .

Let  $G_y = (L', R, E_y)$  denote the bipartite graph associated with  $f_y$ . As implied by the proof of Observation 1, the min-entropy of  $f_y(U_\ell)$  is the number of nodes in  $L'$  that are non-isolated in  $G_y$ . Although each node in  $L'$  is non-isolated in  $G$  (since  $M \subseteq E$ ), edges incident to  $L'$  may disappear when we hardwire  $y$ . We claim that with high probability over  $y$ , plenty of nodes in  $L'$  are still non-isolated in  $G_y$  and hence  $f_y(U_\ell)$  has high min-entropy. For  $i \in \{1, \dots, \ell\}$  let  $(j_i, \text{out}) \in R$  be the neighbor of  $(i, \text{in})$  in  $M$ , and let  $I_{j_i} \times \{\text{in}\}$  be the set of neighbors of  $(j_i, \text{out})$  in  $G$ . Since the  $j_i^{\text{th}}$  output bit of  $f$  depends on the  $i^{\text{th}}$  input bit, there exists a string  $w_i \in \{0, 1\}^{|I_{j_i}|-1}$  such that hardwiring the input bits corresponding to  $I_{j_i} \setminus \{i\}$  to  $w_i$  leaves the edge  $\{(i, \text{in}), (j_i, \text{out})\}$  in place, and in particular ensures that  $(i, \text{in})$  is non-isolated. Since  $M$  is superindependent, the sets  $I_{j_i}$  for  $i \in \{1, \dots, \ell\}$  are pairwise disjoint and in particular, each  $I_{j_i} \setminus \{i\} \subseteq \{\ell + 1, \dots, r\}$ . We assume the bits of  $y$  are indexed starting at  $\ell + 1$ , so for example  $y|_{\{\ell+1\}}$  is the first bit of  $y$ . By the disjointness, we find that the events  $y|_{I_{j_i} \setminus \{i\}} = w_i$  (for  $i \in \{1, \dots, \ell\}$ ) are fully independent over  $y \sim U_{r-\ell}$ . Moreover, each of these events occurs with probability at least  $1/2^{d-1}$  since  $|w_i| \leq d - 1$ . Thus we have

$$\begin{aligned} & \Pr_{y \sim U_{r-\ell}} [f_y(U_\ell) \text{ does not have min-entropy at least } k'] \\ = & \Pr_{y \sim U_{r-\ell}} \left[ \left| \{i \in \{1, \dots, \ell\} : (i, \text{in}) \text{ is non-isolated in } G_y\} \right| < k' \right] \\ \leq & \Pr_{y \sim U_{r-\ell}} \left[ \left| \{i \in \{1, \dots, \ell\} : y|_{I_{j_i} \setminus \{i\}} = w_i\} \right| < k' \right] \\ \leq & e^{-k/8d^2c2^d} \end{aligned}$$

by a standard Chernoff bound.

To summarize, we have shown that every  $(d, c)$ -local source with min-entropy at least  $k/2$  is a uniform convex combination of  $(1, c)$ -local sources, at most  $e^{-k/8d^2c2^d}$  fraction of which do not have min-entropy at least  $k'$ . It now follows from Corollary 1 that  $\text{Ext}$  is a  $(k/2, \epsilon)$ -extractor for  $(d, c)$ -local sources. This finishes the proof of Theorem 7.

## 5 Increasing the Output Length

Combining the results from Section 3 and Section 4 yields an extractor for  $d$ -local sources with output length  $\Omega(k^2/nd^32^d)$ , provided  $d \leq o(\log n)$  and the min-entropy  $k$  is at least  $n^{2/3+\gamma}$ . In this section we show how to improve the output length to  $\Omega(k^2/nd)$ , which is a significant improvement when  $k \geq \Omega(n)$  and  $d$  is large. We present the general method in Section 5.1, and then we apply the general method to obtain Theorem 1 in Section 5.2.

### 5.1 The General Method

We now present our general theorem on increasing the output length of extractors for  $d$ -local sources (Theorem 8 below), which uses the technique of “obtaining an independent seed”. As in [GRS06], the strategy is to take the output of a deterministic extractor and use part of it to sample a set of coordinates of the source, which are then plugged into a seeded extractor, using the other part of the deterministic extractor’s output as the seed. The key property of  $d$ -local sources that enables us to adapt the technique from [GRS06] is that conditioning on any  $p$  bits of the source gives a convex combination of  $d$ -local sources that lose at most  $pd$  in the min-entropy.

A key ingredient (which was not used in [GRS06]) is a fundamental lemma of Nisan and Zuckerman [NZ96], which roughly says that if we sample the coordinates appropriately, then the min-entropy rate of the marginal distribution on those coordinates is almost as high as the min-entropy rate of the whole source.<sup>3</sup> However, the original Nisan-Zuckerman lemma loses a logarithmic factor in the min-entropy rate. We use a strengthened version of the lemma, due to Vadhan [Vad04], which only loses a constant factor.

We use  $\binom{\{1, \dots, n\}}{p}$  to denote the set of subsets of  $\{1, \dots, n\}$  of size  $p$ .

**Definition 2.** We say  $\text{Samp} : \{0, 1\}^s \rightarrow \binom{\{1, \dots, n\}}{p}$  is a  $(\mu, \eta)$ -sampler if for every  $g : \{1, \dots, n\} \rightarrow [0, 1]$  with  $\frac{1}{n} \sum_{j=1}^n g(j) \geq \mu$  it holds that  $\Pr_{\sigma \sim U_s} \left[ \frac{1}{p} \sum_{j \in \text{Samp}(\sigma)} g(j) < \mu/2 \right] \leq \eta$ .

**Lemma 4 ([Vad04]).** There exists a universal constant  $\alpha > 0$  such that the following holds. Suppose  $\text{Samp} : \{0, 1\}^s \rightarrow \binom{\{1, \dots, n\}}{p}$  is a  $(k/2n \log(4n/k), \eta)$ -sampler and  $D$  is a distribution on  $\{0, 1\}^n$  with min-entropy at least  $k$ . Then with probability at least  $1 - \sqrt{\eta + 2^{-\alpha k}}$  over  $\sigma \sim U_s$  it holds that  $D|_{\text{Samp}(\sigma)}$  is  $\sqrt{\eta + 2^{-\alpha k}}$ -close to a distribution with min-entropy at least  $pk/4n$ .

We also need the following lemma from [GRS06], which we state in a slightly nonstandard way for convenience when we apply the lemma.

**Lemma 5 ([GRS06]).** Consider any distribution on  $\{0, 1\}^{s_1} \times \{0, 1\}^{s_2} \times \{0, 1\}^{s_3}$  which is  $\epsilon'$ -close to uniform, and suppose  $\sigma$  is in the support of the marginal distribution on the second coordinate. Then the marginal distribution on the first and third coordinates, conditioned on the second coordinate being  $\sigma$ , is  $(\epsilon'2^{s_2+1})$ -close to uniform.

We now present the general theorem on increasing the output length.

**Theorem 8.** Consider the construction in Figure 1, and let  $\alpha$  be as in Lemma 4. Suppose  $\text{Ext}'$  is a  $(k', \epsilon')$ -extractor for  $d$ -local sources,  $\text{Samp}$  is a  $(k/2n \log(4n/k), \eta)$ -sampler, and  $\text{SExt}$  is a seeded  $(pk/4n, \epsilon'')$ -extractor. Then  $\text{Ext}$  is a  $(k, \epsilon)$ -extractor for  $d$ -local sources, where  $k = k' + pd$  and  $\epsilon = \epsilon'(2^{s+1} + 1) + 2\sqrt{\eta + 2^{-\alpha k}} + \epsilon''$ .

---

<sup>3</sup>Min-entropy rate just means the min-entropy divided by the length of the source.

<p><b>Ingredients:</b></p> <p><math>\text{Ext}' : \{0, 1\}^n \rightarrow \{0, 1\}^{m'}</math></p> <p><math>\text{Ext}'_1 : \{0, 1\}^n \rightarrow \{0, 1\}^s</math> is the first <math>s</math> bits of <math>\text{Ext}'</math></p> <p><math>\text{Ext}'_2 : \{0, 1\}^n \rightarrow \{0, 1\}^{m'-s}</math> is the last <math>m' - s</math> bits of <math>\text{Ext}'</math></p> <p><math>\text{Samp} : \{0, 1\}^s \rightarrow \binom{\{1, \dots, n\}}{p}</math></p> <p><math>\text{SExt} : \{0, 1\}^p \times \{0, 1\}^{m'-s} \rightarrow \{0, 1\}^m</math></p> <p><b>Result:</b></p> <p><math>\text{Ext} : \{0, 1\}^n \rightarrow \{0, 1\}^m</math> defined as <math>\text{Ext}(z) = \text{SExt}(z _{\text{Samp}(\text{Ext}'_1(z))}, \text{Ext}'_2(z))</math></p>
--

Figure 1: Increasing the output length of an extractor for  $d$ -local sources

*Proof.* Consider an arbitrary  $d$ -local sampler  $f : \{0, 1\}^r \rightarrow \{0, 1\}^n$  whose output distribution has min-entropy at least  $k$ , and let  $G = (L, R, E)$  be the associated bipartite graph. Our goal is to show that  $\|\text{Ext}(f(U_r)) - U_m\| \leq \epsilon$ .

Let us call  $\sigma \in \{0, 1\}^s$  *good* if  $f(U_r)|_{\text{Samp}(\sigma)}$  is  $\sqrt{\eta + 2^{-\alpha k}}$ -close to a distribution with min-entropy at least  $pk/4n$ , and *bad* otherwise. For each  $\sigma$  we let  $U_r^{(\sigma)}$  be the uniform distribution over  $w \in \{0, 1\}^r$  such that  $\text{Ext}'_1(f(w)) = \sigma$ .<sup>4</sup>

**Claim 1.** For each good  $\sigma$ ,  $\|\text{Ext}(f(U_r^{(\sigma)})) - U_m\| \leq \epsilon' 2^{s+1} + \sqrt{\eta + 2^{-\alpha k}} + \epsilon''$ .

Assuming Claim 1, we can prove the theorem as follows. Observe that

$$f(U_r) = \sum_{\sigma \in \{0, 1\}^s} \Pr_{w \sim U_r} [\text{Ext}'_1(f(w)) = \sigma] f(U_r^{(\sigma)}).$$

Then using the shorthand  $\epsilon''' = \epsilon' 2^{s+1} + \sqrt{\eta + 2^{-\alpha k}} + \epsilon''$  we have

$$\begin{aligned} \|\text{Ext}(f(U_r)) - U_m\| &\leq \epsilon''' + \Pr_{w \sim U_r} \left[ \|\text{Ext}(f(U_r^{(\sigma)})) - U_m\| > \epsilon''' \text{ where } \sigma = \text{Ext}'_1(f(w)) \right] \\ &\leq \epsilon''' + \Pr_{w \sim U_r} [\text{Ext}'_1(f(w)) \text{ is bad}] \\ &\leq \epsilon''' + \epsilon' + \Pr_{\sigma \sim U_s} [\sigma \text{ is bad}] \\ &\leq \epsilon''' + \epsilon' + \sqrt{\eta + 2^{-\alpha k}} \\ &= \epsilon \end{aligned}$$

where the first line follows by Lemma 1, the second line follows by Claim 1, the third line follows by  $\|\text{Ext}'_1(f(U_r)) - U_s\| \leq \epsilon'$  (since  $f(U_r)$  is a  $d$ -local source with min-entropy at least  $k \geq k'$ ), and the fourth line follows by Lemma 4.

It remains to prove Claim 1. Consider an arbitrary fixed good  $\sigma \in \{0, 1\}^s$ , and without loss of generality assume the nodes in  $L$  adjacent to  $\text{Samp}(\sigma) \times \{\text{out}\}$  are  $\{r - \ell + 1, \dots, r\} \times \{\text{in}\}$  for some  $\ell \leq pd$ . For each string  $y \in \{0, 1\}^\ell$ , define  $f_y : \{0, 1\}^{r-\ell} \rightarrow \{0, 1\}^n$  as  $f_y(x) = f(x, y)$  (hardwiring the last  $\ell$  bits to  $y$ ). Then each  $f_y(U_{r-\ell})$  is a  $d$ -local source with min-entropy at least  $k - \ell \geq k'$  (see the proof of Lemma 2), and this is the key point that enables us to use the technique of [GRS06]. Thus,  $\|\text{Ext}'(f_y(U_{r-\ell})) - U_{m'}\| \leq \epsilon'$ . Now consider the joint distribution

$$\left( U_r|_{\{r-\ell+1, \dots, r\}}, \text{Ext}'_1(f(U_r)), \text{Ext}'_2(f(U_r)) \right).$$

<sup>4</sup>Formally, we only consider  $\sigma$ 's in the support of  $\text{Ext}'_1(f(U_r))$ .

That is, sample  $(x, y) \sim U_r$  and output  $y$  along with both parts of  $\text{Ext}'(f(x, y))$ . We have just argued that conditioned on the first coordinate of this distribution being any particular  $y \in \{0, 1\}^\ell$ , the marginal distribution of the other two coordinates is  $\epsilon'$ -close to uniform. Thus the entire distribution is  $\epsilon'$ -close to uniform. By Lemma 5 (with  $s_1 = \ell$ ,  $s_2 = s$ , and  $s_3 = m' - s$ ), the joint distribution

$$\left( U_r^{(\sigma)}|_{\{r-\ell+1, \dots, r\}}, \text{Ext}'_2(f(U_r^{(\sigma)})) \right)$$

is  $(\epsilon'2^{s+1})$ -close to the uniform distribution  $(U_\ell, U_{m'-s})$  where  $U_\ell$  and  $U_{m'-s}$  are independent. Let us define  $f^{(\sigma)} : \{0, 1\}^\ell \rightarrow \{0, 1\}^p$  by  $f^{(\sigma)}(y) = f(x, y)|_{\text{Samp}(\sigma)}$  for any  $x \in \{0, 1\}^{r-\ell}$  (this value does not depend on  $x$  since nodes in  $\text{Samp}(\sigma) \times \{\text{out}\}$  are only adjacent to nodes in  $\{r-\ell+1, \dots, r\} \times \{\text{in}\}$ ). Then we have

$$\text{Ext}(f(U_r^{(\sigma)})) = \text{SExt}\left(f^{(\sigma)}(U_r^{(\sigma)}|_{\{r-\ell+1, \dots, r\}}), \text{Ext}'_2(f(U_r^{(\sigma)}))\right)$$

and thus

$$\|\text{Ext}(f(U_r^{(\sigma)})) - \text{SExt}(f^{(\sigma)}(U_\ell), U_{m'-s})\| \leq \epsilon'2^{s+1}. \quad (1)$$

Letting  $D$  denote a distribution on  $\{0, 1\}^p$  with min-entropy at least  $pk/4n$  that  $f(U_r)|_{\text{Samp}(\sigma)} = f^{(\sigma)}(U_\ell)$  is  $\sqrt{\eta + 2^{-\alpha k}}$ -close to (such a  $D$  exists since  $\sigma$  is good), we have

$$\|\text{SExt}(f^{(\sigma)}(U_\ell), U_{m'-s}) - \text{SExt}(D, U_{m'-s})\| \leq \sqrt{\eta + 2^{-\alpha k}}. \quad (2)$$

Since  $\text{SExt}$  is a seeded  $(pk/4n, \epsilon'')$ -extractor, we have

$$\|\text{SExt}(D, U_{m'-s}) - U_m\| \leq \epsilon''. \quad (3)$$

Combining Inequality (1), Inequality (2), and Inequality (3) yields Claim 1. This finishes the proof of Theorem 8.  $\square$

## 5.2 Applying Theorem 8

In order to apply Theorem 8, we need explicit constructions of  $\text{Ext}'$ ,  $\text{Samp}$ , and  $\text{SExt}$ . An appropriate construction of  $\text{Samp}$  is given by the following lemma.

**Lemma 6 ([NZ96]).** *There exists an explicit  $(\mu, \eta)$ -sampler  $\text{Samp} : \{0, 1\}^s \rightarrow \binom{\{1, \dots, n\}}{p}$  with  $s = 4 \log n \cdot \log \frac{1}{\eta}$ , provided  $\mu p \geq 64 \log \frac{1}{\eta}$  and  $\eta < 1/16$ .*

The interesting thing about samplers as defined in Definition 2 is that they produce a set of fixed size. (Typically, samplers produce either a multiset of fixed size or a set of random size, and the latter is sufficient for the argument in [GRS06].) Nisan and Zuckerman [NZ96] proved Lemma 6 by partitioning the  $n$  coordinates into  $p$  blocks, picking one coordinate from each block in an  $O(\log \frac{1}{\eta})$ -wise independent way, and using the concentration bounds of [SSS95, BR94].<sup>5</sup> Vadhan [Vad04] also constructed a sampler that produces a set of fixed size, and with better seed length for a certain range of parameters. However, his seed length is actually not good enough for our range of parameters.

As for the seeded extractor  $\text{SExt}$ , plenty of known constructions are good enough for our purpose. For example, we can use the following construction, due to Raz, Reingold, and Vadhan.

<sup>5</sup>Actually, Nisan and Zuckerman proved a version with slightly different constants and where the sampler only needs to work for boolean functions  $g$ , but the proof goes through to yield Lemma 6.

**Theorem 9 ([RRV99]).** *There exists an explicit seeded  $(k, \epsilon)$ -extractor  $\text{SExt} : \{0, 1\}^n \times \{0, 1\}^t \rightarrow \{0, 1\}^m$  with  $t = O((\log^2 n + \log \frac{1}{\epsilon}) \cdot \log k)$  and  $m = k$ .*

At last, we can prove Theorem 1.

*Proof of Theorem 1.* Assume  $k \geq n^{2/3+\gamma}$  and  $d \leq \beta \log n$  for small enough constant  $\beta > 0$ . Then for some constant  $\beta' > 0$  to be specified shortly, define the following parameters.

- $\epsilon' = 2^{-n^{\beta'}}$
- $k' = k/2$
- $m' = (k')^2 / 8nd^3 2^d$
- $p = k/2d$
- $\mu = k/2n \log(4n/k)$
- $\eta = 2^{-n^{\beta'/2}}$
- $s = 4 \log n \cdot \log \frac{1}{\eta}$
- $\epsilon'' = 2^{-n^{1/4}}$
- $m = pk/4n$
- $t = m' - s$

As shown in the discussion after the statement of Theorem 7, combining Theorem 7 with Theorem 5 yields an explicit  $(k', \epsilon')$ -extractor  $\text{Ext}' : \{0, 1\}^n \rightarrow \{0, 1\}^{m'}$  for  $d$ -local sources, provided  $\beta'$  is small enough. By Lemma 6 there exists an explicit  $(\mu, \eta)$ -sampler  $\text{Samp} : \{0, 1\}^s \rightarrow \binom{\{1, \dots, n\}}{p}$ . Since  $t \geq \omega((\log^2 p + \log \frac{1}{\epsilon'}) \cdot \log m)$ , by Theorem 9 there exists an explicit seeded  $(m, \epsilon'')$ -extractor  $\text{SExt} : \{0, 1\}^p \times \{0, 1\}^t \rightarrow \{0, 1\}^m$ . Thus by Theorem 8,  $\text{Ext}$  is a  $(k, \epsilon)$ -extractor for  $d$ -local sources, where  $\epsilon = \epsilon'(2^{s+1} + 1) + 2\sqrt{\eta + 2^{-\alpha k}} + \epsilon'' \leq 2^{-n^\beta}$  provided  $\beta$  is small enough.  $\square$

## 6 Improved Lower Bounds for Sampling Input-Output Pairs

For this section, we define a  $(d, c, k)$ -local sampler to be a  $(d, c)$ -local sampler with at least  $k$  non-isolated nodes on the left side of its associated bipartite graph (that is, it makes nontrivial use of at least  $k$  random bits). We say a distribution on  $\{0, 1\}^n$  is a  $(d, c, k)$ -local source if it equals  $f(U_r)$  for some  $(d, c, k)$ -local sampler  $f$  (with any input length  $r$ ). Note that a  $(d, c, k)$ -local source might not have min-entropy at least  $k$ .

**Theorem 10.** *Suppose  $\text{Ext} : \{0, 1\}^n \rightarrow \{0, 1\}$  is a  $(0, \epsilon)$ -extractor for  $(d, 8d, n/4)$ -local sources, where  $d < n/8$ . Then for every  $d$ -local source  $D$  on  $\{0, 1\}^{n+1}$  we have  $\|D - (U_n, \text{Ext}(U_n))\| \geq 1/2 - \epsilon - 2^{-n/2}$ .*

It might seem suspicious that we are assuming Ext is a  $(0, \epsilon)$ -extractor. We are not, in fact, extracting from sources with 0 min-entropy — it is possible to derive a lower bound on the min-entropy of any  $(d, 8d, n/4)$ -local source.<sup>6</sup> The point is that for Theorem 10, we do not care about the min-entropy, only the number of non-isolated input nodes. Before proving Theorem 10, we show how it implies Theorem 3.

*Proof of Theorem 3.* In the proof of Theorem 7, we implicitly showed that for all  $n, k, d, c, \epsilon'$ , every  $(k', \epsilon')$ -extractor for  $(1, c)$ -local sources is also a  $(k, \epsilon)$ -extractor for  $(d, c)$ -local sources where  $k' = k/d^2 c 2^d$  and  $\epsilon = \epsilon' + e^{-k'/4}$  (by replacing  $k/2$  with  $k$  in the proof). The only property of having min-entropy at least  $k$  we used in that proof was that the sampler must make nontrivial use of at least  $k$  random bits; thus we can conclude that the extractor is a  $(0, \epsilon)$ -extractor for  $(d, c, k)$ -local sources.

Assume  $d \leq \beta \log n$  for some small enough constant  $\beta > 0$ . Set  $c = 8d$  and  $k = n/4$  and  $k' = k/d^2 c 2^d = n/32d^3 2^d \geq n^{1/2}$ . Using  $\gamma = 1/2$  in Theorem 5, there exists an explicit  $(k', \epsilon')$ -extractor for  $(1, c)$ -local sources with output length 1, where  $\epsilon' = 2^{-n^{\Omega(1)}}$  (since  $k' \geq \log^{\omega(1)} n$  and  $(k')^{1/2} \geq c$  and  $k' - o(k') \geq 1$ ). By the observation in the previous paragraph, this function is a  $(0, \epsilon)$ -extractor for  $(d, 8d, n/4)$ -local sources with error  $\epsilon = 2^{-n^{\Omega(1)}}$ . Theorem 3 follows immediately from this and Theorem 10.  $\square$

*Proof of Theorem 10.* Consider an arbitrary  $d$ -local sampler  $f : \{0, 1\}^r \rightarrow \{0, 1\}^{n+1}$ , and let  $G = (L, R, E)$  be the associated bipartite graph. Since  $|E| \leq (n+1)d$ , there are at most  $(n+1)/8$  nodes in  $L$  with degree greater than  $8d$ . Also, at most  $d \leq (n-1)/8$  nodes in  $L$  are adjacent to  $(n+1, \text{out})$ . Without loss of generality, the nodes in  $L$  that either have degree greater than  $8d$  or are adjacent to  $(n+1, \text{out})$  are  $\{r - \ell + 1, \dots, r\} \times \{\text{in}\}$  for some  $\ell \leq (n+1)/8 + (n-1)/8 = n/4$ . For each string  $y \in \{0, 1\}^\ell$ , define  $f_y : \{0, 1\}^{r-\ell} \rightarrow \{0, 1\}^{n+1}$  as  $f_y(x) = f(x, y)$  (hardwiring the last  $\ell$  bits to  $y$ ) and let  $G_y = (L', R, E_y)$  be the associated bipartite graph, where  $L' = \{1, \dots, r - \ell\} \times \{\text{in}\}$ . Observe that  $f(U_r) = \sum_{y \in \{0, 1\}^\ell} \frac{1}{2^\ell} f_y(U_{r-\ell})$ . We define the tests

$$T_1 = \left\{ z \in \{0, 1\}^{n+1} : \exists x \in \{0, 1\}^{r-\ell}, y \in \{0, 1\}^\ell \text{ such that } f(x, y) = z \text{ and } \left| \{i \in \{1, \dots, r - \ell\} : (i, \text{in}) \text{ is non-isolated in } G_y\} \right| < n/4 \right\}$$

and

$$T_2 = \left\{ z \in \{0, 1\}^{n+1} : \text{Ext}(z|_{\{1, \dots, n\}}) \neq z|_{\{n+1\}} \right\}$$

(in other words, the support of  $(U_n, \text{Ext}(U_n))$  is the complement of  $T_2$ ). Finally, we define the test  $T = T_1 \cup T_2$ .

**Claim 2.**  $\Pr_{f(U_r)}[T] \geq 1/2 - \epsilon$ .

**Claim 3.**  $\Pr_{(U_n, \text{Ext}(U_n))}[T] \leq 2^{-n/2}$ .

Combining the two claims, we have  $|\Pr_{f(U_r)}[T] - \Pr_{(U_n, \text{Ext}(U_n))}[T]| \geq 1/2 - \epsilon - 2^{-n/2}$ , thus witnessing that  $\|f(U_r) - (U_n, \text{Ext}(U_n))\| \geq 1/2 - \epsilon - 2^{-n/2}$ .

---

<sup>6</sup>Specifically, a combinatorial argument shows that the source must have many bits that are fully independent of each other and that each have probability  $\geq 1/2^d$  for both outcomes 0 and 1. A lower bound on the min-entropy can be derived from this fact.

*Proof of Claim 2.* It suffices to show that for each  $y \in \{0, 1\}^\ell$ ,  $\Pr_{f_y(U_{r-\ell})}[T] \geq 1/2 - \epsilon$ . If  $y$  is such that  $|\{i \in \{1, \dots, r - \ell\} : (i, \text{in}) \text{ is non-isolated in } G_y\}| < n/4$  then of course  $\Pr_{f_y(U_{r-\ell})}[T_1] = 1$ . Otherwise,  $f_y(U_{r-\ell})$  is a  $(d, 8d, n/4)$ -source on  $\{0, 1\}^{n+1}$ . Note that  $(n + 1, \text{out})$  is isolated in  $G_y$ ; we define  $b_y \in \{0, 1\}$  to be the fixed value of the  $(n + 1)^{\text{st}}$  output bit of  $f_y$ , and we define  $f'_y : \{0, 1\}^{r-\ell} \rightarrow \{0, 1\}^n$  to be the first  $n$  output bits of  $f_y$ . Since  $f'_y(U_{r-\ell})$  is a  $(d, 8d, n/4)$ -source on  $\{0, 1\}^n$ , we have  $\|\text{Ext}(f'_y(U_{r-\ell})) - U_1\| \leq \epsilon$  and thus  $\Pr_{b \sim \text{Ext}(f'_y(U_{r-\ell}))}[b \neq b_y] \geq 1/2 - \epsilon$ . In other words,  $\Pr_{f_y(U_{r-\ell})}[T_2] \geq 1/2 - \epsilon$ . This finishes the proof of Claim 2.  $\square$

*Proof of Claim 3.* By definition,  $\Pr_{(U_n, \text{Ext}(U_n))}[T_2] = 0$ . Note that  $|T_1| \leq 2^{n/2}$  since each string in  $T_1$  can be described by a string of length at most  $\ell + n/4 \leq n/2$ , namely an appropriate value of  $y$  along with the bits of  $x$  such that the corresponding nodes in  $L'$  are non-isolated in  $G_y$ . Since  $(U_n, \text{Ext}(U_n))$  is uniform over a set of size  $2^n$ , we get  $\Pr_{(U_n, \text{Ext}(U_n))}[T_1] \leq 2^{n/2}/2^n = 2^{-n/2}$ . This finishes the proof of Claim 3.  $\square$

This finishes the proof of Theorem 10.  $\square$

## 7 Improved Extractors for Low-Weight Affine Sources

We now describe the proof of Theorem 5. To do this, we need a construction of linear strong seeded extractors with good seed length, which we present in Section 7.1. Then in Section 7.2 we derive Theorem 5.

### 7.1 A Linear Strong Seeded Extractor with Seed Length $\log n + O(\log k)$

Our goal in this section is to prove the following theorem.<sup>7</sup>

**Theorem 11.** *There exists a constant  $c$  such that for all  $k \geq c \log^2 n$  there exists an explicit linear strong seeded  $(k, 1/4)$ -extractor  $\text{Ext} : \{0, 1\}^n \times \{0, 1\}^t \rightarrow \{0, 1\}^m$  with  $t = \log n + c \log k$  and  $m = k^{1/4}$ .*

It is very important to us that the seed length here has  $\log n$  and not  $O(\log n)$ . If instead we use an extractor with  $c \log n$  in the seed length, then in Theorem 5 we would only be able to get an extractor for the class of weight- $k^{(1/c)-\gamma}$  affine sources as opposed to weight- $k^{1-\gamma}$  affine sources.

We also note that without the linearity property, such an extractor is explicitly constructed and stated in [GUV09, Theorem 5.12]. We construct such an extractor with the linearity property by using a construction from [GUV09] and then bootstrapping it with another known construction. To do this, we first define and construct objects called linear strong condensers.

**Definition 3.** *We say  $C : \{0, 1\}^n \times \{0, 1\}^t \rightarrow \{0, 1\}^m$  is a strong  $k \rightarrow_\epsilon k'$  condenser if for every distribution  $D$  on  $\{0, 1\}^n$  with min-entropy at least  $k$ ,*

$$\Pr_{y \sim U_t} [C(D, y) \text{ is } \epsilon\text{-close to a distribution with min-entropy at least } k'] \geq 1 - \epsilon.$$

---

<sup>7</sup>We note that Theorem 11 can be generalized by setting the parameters appropriately in our argument. For general error  $\epsilon$ , we can get seed length  $\log n + O(\log(k/\epsilon))$ , and the output length can be improved to  $k^{1-\alpha}$  for any constant  $\alpha > 0$  at the expense of increasing the lower bound on  $k$ . However, we only prove the version we need for Theorem 5.

Recall that we say  $C$  is *linear* if for every  $y \in \{0, 1\}^t$ , the function  $C(\cdot, y) : \{0, 1\}^n \rightarrow \{0, 1\}^m$  is linear over  $\mathbb{F}_2$ .

Our construction of a linear strong condenser is the same as one of the constructions in [GUV09], which in turn is based on an idea from [GR08b]. However, we need to argue about its linearity as well as the parameters, so we state the construction and result of [GUV09] here. Consider a finite field  $\mathbb{F}_q$  for some  $q = 2^t$ . Let  $\zeta$  be a generator of the multiplicative group  $\mathbb{F}_q^*$ . Then the function  $C : \mathbb{F}_q^{n'} \times \mathbb{F}_q \rightarrow \mathbb{F}_q^{m'}$  is as follows.

Given  $f = (f_0, \dots, f_{n'-1}) \in \mathbb{F}_q^{n'}$ , we interpret it as a polynomial  $f : \mathbb{F}_q \rightarrow \mathbb{F}_q$  such that  $f : y \mapsto \sum_{0 \leq i < n'} f_i y^i$ . We now describe  $C$  as  $C : (f, y) \mapsto (f(y), f(\zeta y), \dots, f(\zeta^{m'-1} y))$ .

**Observation 2.** *For all  $y \in \mathbb{F}_q$ , the function  $C(\cdot, y) : f \mapsto C(f, y)$  is  $\mathbb{F}_q$ -linear.*

**Observation 3.** *For  $q = 2^t$ , there is an isomorphism between  $(\mathbb{F}_q, +)$  and  $(\mathbb{F}_2^t, \oplus)$ . Further, this isomorphism is computable in time polynomial in  $t$ .*

Thus we can interpret  $C$  as a function  $C : \{0, 1\}^n \times \{0, 1\}^t \rightarrow \{0, 1\}^m$  where  $n = n' \cdot t$  and  $m = m' \cdot t$ . Further,  $C$  is  $\mathbb{F}_2$ -linear, and it is polynomial time computable since a generator of  $\mathbb{F}_q^*$  can be computed in time polynomial in  $t$  [Sho88]. The fact that  $C$  is a strong condenser follows from [GUV09, Theorem 7.2].

**Theorem 12 ([GUV09]).** *For every  $\ell \leq n$  such that  $2^\ell$  is an integer, and for every  $\alpha, \epsilon > 0$ , the function  $C : \{0, 1\}^n \times \{0, 1\}^t \rightarrow \{0, 1\}^m$  as defined above is a*

$$\text{strong } (1 + 1/\alpha)\ell d + \log(1/\epsilon) \rightarrow_{\sqrt{3\epsilon}} \ell d - 2 \text{ condenser}$$

*with  $t \leq (1 + 1/\alpha)d$  and  $m \leq (1 + 1/\alpha)\ell d$  where  $d = \lceil \alpha \log(4n\ell/\epsilon) \rceil$ , provided  $\ell d \geq \log(1/\epsilon)$ .*

The following important corollary follows by setting parameters correctly in the result. Assume  $k \geq c \log^2 n$  for some large constant  $c > 0$ , and set the parameters as follows.

- $\epsilon = 1/2^8$
- $\alpha = (\log k)/(2^8 \cdot \log n)$
- $\ell = k/(2^8 \cdot \log n)$

This implies the following.

- $d = \left\lceil \frac{(\log k)(\log n + \log k - \log \log n + O(1))}{2^8 \cdot \log n} \right\rceil = \frac{(\log k)(\log n + \log k - \log \log n)(1 + o(1))}{2^8 \cdot \log n}$
- $t \leq (\log n + \log k - \log \log n)(1 + (2 \log k / \log n)) \leq \log n + 5 \log k$
- $m \leq \frac{k}{2^8 \cdot \log n} \cdot (\log n + 5 \log k) \leq k$
- $\ell d \geq (k \log k)/(2^{16} \cdot \log n) \geq k^{1/2} + 2$
- $k \geq (1 + 1/\alpha)\ell d + \log(1/\epsilon)$

Hence, we now get the following corollary.



**Corollary 4.** *There exists a constant  $c$  such that for all  $k \geq c \log^2 n$  there exists an explicit linear strong  $k \rightarrow_{1/8} k^{1/2}$  condenser  $C : \{0, 1\}^n \times \{0, 1\}^t \rightarrow \{0, 1\}^m$  with  $t = \log n + 5 \log k$  and  $m = k$ .*

We now recall that the strong seeded extractors in [Tre01, RRV02] are also linear.

**Theorem 13 ([Tre01]).** *There exists an explicit linear strong seeded  $(n^{1/2}, 1/8)$ -extractor  $\text{Ext} : \{0, 1\}^n \times \{0, 1\}^t \rightarrow \{0, 1\}^m$  with  $t = O(\log n)$  and  $m = n^{1/4}$ .*

Theorem 11 follows from Corollary 4 and Theorem 13.

## 7.2 Proof of Theorem 5

In this section we prove Theorem 5. As we have said before, Rao [Rao09b] proves the same kind of theorem except it is weaker in the upper bound on the weight allowed for the affine sources. Our extractor construction uses the same steps as [Rao09b], except the components used in our construction are tailor-made for our purposes thus helping us achieve better parameters. Throughout this section, all references to particular theorems in [Rao09b] actually refer to the ECC version of the paper (technical report TR08-015). Also, throughout this section we let  $c$  be the constant from Theorem 11.

In order to describe the better extractors, we first recall the following linear error-correcting code construction (BCH code) [Sud].

**Theorem 14.** *For every  $d < n$  there exists an explicit parity check function  $P : \mathbb{F}_2^n \rightarrow \mathbb{F}_2^m$  for a linear code with distance greater than  $d$ , such that  $m = O(d \log n)$ .*

We now recall the following claim from [Rao09b, Lemma 6.1].

**Claim 4.** *Let  $P : \mathbb{F}_2^n \rightarrow \mathbb{F}_2^m$  be a parity check function for a linear code with distance greater than  $d$ . Let  $D$  be any weight- $w$  affine source with min-entropy at least  $d/w$ . Then  $P(D)$  is an affine source with min-entropy at least  $d/w$ .*

Combining Theorem 14 and Claim 4 (using  $d = k^{1-\gamma/2}$ ), we get the following.

**Lemma 7.** *For every constant  $\gamma > 0$  and all  $k$  there exists an explicit linear function  $P : \mathbb{F}_2^n \rightarrow \mathbb{F}_2^m$  with  $m = O(k^{1-\gamma/2} \cdot \log n)$  such that if  $D$  is a weight- $k^{1-\gamma}$  affine source with min-entropy at least  $k$ , then  $P(D)$  is an affine source with min-entropy at least  $k^{\gamma/2}$ .*

Now let  $\gamma$  and  $k$  be as in Theorem 5. Let  $m_0 = O(k^{1-\gamma/2} \cdot \log n)$  be the output length from Lemma 7. Let  $\text{Ext}_1 : \{0, 1\}^{m_0} \times \{0, 1\}^{t_1} \rightarrow \{0, 1\}^{m_1}$  be the linear strong seeded extractor from Theorem 11 set up to work for min-entropy  $k^{\gamma/4c}$  (which is less than  $k^{\gamma/2}$  and is at least  $c \log^2 m_0$  assuming  $k \geq \log^{10c/\gamma} n$ ). Thus we have  $t_1 = \log m_0 + (\gamma/4) \log k$  and  $m_1 = k^{\gamma/16c}$ . In Figure 2 we present the routine LOW-CONVERT from [Rao09b]. The following lemma was proven in [Rao09b, Lemma 6.3]. We note that for this, the error of  $\text{Ext}_1$  only needs to be  $< 1/2$ .

**Definition 4.** *A distribution  $D$  on  $\{0, 1\}^{\ell \times \ell}$  is said to be an affine somewhere random source if  $D$  is an affine source and for some  $1 \leq i \leq \ell$ , the  $i^{\text{th}}$  row of  $D$  is uniformly random.*

**Lemma 8.** *For every constant  $\gamma > 0$ , if  $D$  is a weight- $k^{1-\gamma}$  affine source with min-entropy at least  $k \geq \log^{10c/\gamma} n$ , then  $\text{LC}(D)$  is an affine somewhere random source of size  $2^{t_1} \times m_1$ .*

LOW-CONVERT( $D$ )

Input:  $x \in \{0, 1\}^n$

Output:  $z \in \{0, 1\}^{2^{t_1} \times m_1}$

Subroutines used:  $P : \{0, 1\}^n \rightarrow \{0, 1\}^{m_0}$  from Lemma 7, and  $\text{Ext}_1 : \{0, 1\}^{m_0} \times \{0, 1\}^{t_1} \rightarrow \{0, 1\}^{m_1}$  from Theorem 11. Here  $m_0 = O(k^{1-\gamma/2} \cdot \log n)$ ,  $t_1 = \log m_0 + (\gamma/4) \log k$ , and  $m_1 = k^{\gamma/16c}$ .

For  $1 \leq i \leq 2^{t_1}$ , the  $i^{\text{th}}$  row of the output is defined by  $LC(x)_i = \text{Ext}_1(P(x), i)$ .

Figure 2: LOW-CONVERT

Note that the number of rows in the output of  $LC$  is  $2^{t_1} = m_0 \cdot k^{\gamma/4} = O(k^{1-\gamma/4} \cdot \log n) \leq k^{1-\gamma/8}$ . At this stage, we also point out how Theorem 11's optimized dependence of the seed length on the length of the source is crucial for the construction. For the rest of the argument to go through, we require the number of rows in  $LC(x)$  (namely  $2^{t_1}$ ) to be smaller than  $k$ . If we used a linear strong seeded extractor for which  $t_1 \geq c' \log m_0$  then this would force  $m_0 < k^{1/c'}$ . However, our use of Theorem 14 and Claim 4 requires  $m_0 > w$ , which would imply that we need  $w < k^{1/c'}$ . Instead, using our optimized extractor from Theorem 11, we are able to handle any weight  $w \leq k^{1-\gamma}$ .

In order to define the next routine, we recall an extractor construction from [RRV02] for a particular setting of parameters.

**Theorem 15 ([RRV02]).** *There is an explicit linear strong seeded  $(k, \epsilon)$ -extractor  $\text{Ext}_2 : \{0, 1\}^n \times \{0, 1\}^t \rightarrow \{0, 1\}^m$  with  $t = O(\log^3(n/\epsilon))$  and  $m = k - O(\log^3(n/\epsilon))$ .*

We set up  $\text{Ext}_2$  to work for min-entropy  $k - 2^{t_1} \cdot m_1 \geq k - k^{1-\gamma/8+\gamma/16c} = k - o(k)$  and seed length  $t_2 = m_1$  and thus we get output length  $m_2 = k - o(k)$  with error  $2^{-k^{\Omega(1)}}$ . In Figure 3 we present the routine AFFINE-CONVERT from [Rao09b]. The following lemma was proven in [Rao09b, Theorem 6.5].

**Lemma 9.** *For every constant  $\gamma > 0$ , if  $D$  is a weight- $k^{1-\gamma}$  affine source with min-entropy at least  $k \geq \log^{10c/\gamma} n$ , then  $AC(D)$  is  $2^{-k^{\Omega(1)}}$ -close to a convex combination of affine somewhere random sources of size  $2^{t_1} \times m_2$ .*

Lemma 9 says that the output of  $AC(D)$  is close to a convex combination of affine somewhere random sources. Since the length of each row is  $m_2$  and the number of rows is  $2^{t_1} \leq k^{1-\gamma/8} \leq m_2^{1-\gamma/9} \ll m_2$ , we can apply the routine AFFINE-SREXT from [Rao09b]. The following lemma was proven in [Rao09b, Theorem 5.1].<sup>8</sup>

**Lemma 10.** *For every constant  $\alpha > 0$  there exists an explicit function  $A : \{0, 1\}^{k^{1-\alpha} \times k} \rightarrow \{0, 1\}^m$  such that if  $D$  is an affine somewhere random source of size  $k^{1-\alpha} \times k$ , then  $\|A(D) - U_m\| \leq 2^{-k^{\Omega(1)}}$  where  $m = k - o(k)$ .*

Theorem 5 follows from Lemma 9 and Lemma 10.

<sup>8</sup>We note that [Rao09b, Theorem 5.1] discusses affine somewhere random sources of size  $k^{0.7} \times k$ . However, it is straightforward to see that the result just requires the number of rows in the affine somewhere random source to be polynomially smaller than the length of each row.

### AFFINE-CONVERT( $D$ )

Input:  $x \in \{0, 1\}^n$

Output:  $z \in \{0, 1\}^{2^{t_1} \times m_2}$

Subroutines used:  $LC : \{0, 1\}^n \rightarrow \{0, 1\}^{2^{t_1} \times m_1}$  from Lemma 8, and  $\text{Ext}_2 : \{0, 1\}^n \times \{0, 1\}^{t_2} \rightarrow \{0, 1\}^{m_2}$  from Theorem 15. Here  $t_2 = m_1$  and  $m_2 = k - o(k)$ .

For  $1 \leq i \leq 2^{t_1}$ , the  $i^{\text{th}}$  row of the output is defined by  $AC(x)_i = \text{Ext}_2(x, LC(x)_i)$ .

Figure 3: AFFINE-CONVERT

## 8 Open Problems

One open problem is to quantitatively improve our results and those of Viola [Vio11]. This may require more sophisticated tools for understanding the min-entropy of the output distribution of a local sampler.

The key new technique introduced in this paper is to show that certain sources are close to convex combinations of low-weight affine sources, and then apply the extractor of [Rao09b]. This technique is very powerful; Viola [Vio11] has shown that it already encompasses sources samplable by polynomial-size constant-depth circuits. What other classes of sources can this technique handle?

In this paper, we have considered samplers where each output bit only depends on a small number of input bits. What about samplers where each input bit only influences a small number of output bits?

## Acknowledgments

We thank Zeev Dvir, Omer Reingold, Salil Vadhan, and Avi Wigderson for helpful email exchanges. In particular, A.D. would like to thank Omer and Salil for answering his innumerable queries about extractors and Salil for suggesting the use of GUV condensers. We had helpful conversations about this work with Urmila Mahadev, Anup Rao, Luca Trevisan, Gregory Valiant, and Emanuele Viola. We also thank anonymous reviewers for helpful comments and suggestions.

## References

- [ABR11] Benny Applebaum, Andrej Bogdanov, and Alon Rosen. A dichotomy for local small-bias generators. Technical Report TR11-126, Electronic Colloquium on Computational Complexity, 2011.
- [AIK06] Benny Applebaum, Yuval Ishai, and Eyal Kushilevitz. Cryptography in  $\text{NC}^0$ . *SIAM Journal on Computing*, 36(4):845–888, 2006.
- [AIK08] Benny Applebaum, Yuval Ishai, and Eyal Kushilevitz. On pseudorandom generators with linear stretch in  $\text{NC}^0$ . *Computational Complexity*, 17(1):38–69, 2008.

- [App11] Benny Applebaum. Pseudorandom generators with long stretch and low locality from random local one-way functions. Technical Report TR11-007, Electronic Colloquium on Computational Complexity, 2011.
- [ASW09] Sanjeev Arora, David Steurer, and Avi Wigderson. Towards a study of low-complexity graphs. In *Proceedings of the 36th International Colloquium on Automata, Languages and Programming*, pages 119–131, 2009.
- [BIW06] Boaz Barak, Russell Impagliazzo, and Avi Wigderson. Extracting randomness using few independent sources. *SIAM Journal on Computing*, 36(4):1095–1118, 2006.
- [BKS<sup>+</sup>10] Boaz Barak, Guy Kindler, Ronen Shaltiel, Benny Sudakov, and Avi Wigderson. Simulating independence: New constructions of condensers, Ramsey graphs, dispersers, and extractors. *Journal of the ACM*, 57(4), 2010.
- [Bou05] Jean Bourgain. More on the sum-product phenomenon in prime fields and its applications. *International Journal of Number Theory*, 1:1–32, 2005.
- [Bou07] Jean Bourgain. On the construction of affine-source extractors. *Geometric and Functional Analysis*, 1:33–57, 2007.
- [BQ09] Andrej Bogdanov and Youming Qiao. On the security of Goldreich’s one-way function. In *Proceedings of the 13th International Workshop on Randomization and Computation*, pages 392–405, 2009.
- [BR94] Mihir Bellare and John Rompel. Randomness-efficient oblivious sampling. In *Proceedings of the 35th IEEE Symposium on Foundations of Computer Science*, pages 276–287, 1994.
- [BRSW06] Boaz Barak, Anup Rao, Ronen Shaltiel, and Avi Wigderson. 2-source dispersers for sub-polynomial entropy and Ramsey graphs beating the Frankl-Wilson construction. In *Proceedings of the 38th ACM Symposium on Theory of Computing*, pages 671–680, 2006.
- [BSG11] Eli Ben-Sasson and Ariel Gabizon. Extractors for polynomials sources over constant-size fields of small characteristic. Technical Report TR11-129, Electronic Colloquium on Computational Complexity, 2011.
- [CEMT09] James Cook, Omid Etesami, Rachel Miller, and Luca Trevisan. Goldreich’s one-way function candidate and myopic backtracking algorithms. In *Proceedings of the 6th Theory of Cryptography Conference*, pages 521–538, 2009.
- [CFG<sup>+</sup>85] Benny Chor, Joel Friedman, Oded Goldreich, Johan Håstad, Steven Rudich, and Roman Smolensky. The bit extraction problem or  $t$ -resilient functions. In *Proceedings of the 26th IEEE Symposium on Foundations of Computer Science*, pages 396–407, 1985.
- [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.

- [CM01] Mary Cryan and Peter Bro Miltersen. On pseudorandom generators in  $NC^0$ . In *Proceedings of the 26th International Symposium on Mathematical Foundations of Computer Science*, pages 272–284, 2001.
- [DEOR04] Yevgeniy Dodis, Ariel Elbaz, Roberto Oliveira, and Ran Raz. Improved randomness extraction from two independent sources. In *Proceedings of the 8th International Workshop on Randomization and Computation*, pages 334–344, 2004.
- [DG10] Matt DeVos and Ariel Gabizon. Simple affine extractors using dimension expansion. In *Proceedings of the 25th IEEE Conference on Computational Complexity*, pages 50–57, 2010.
- [DGW09] Zeev Dvir, Ariel Gabizon, and Avi Wigderson. Extractors and rank extractors for polynomial sources. *Computational Complexity*, 18(1):1–58, 2009.
- [DM04] Stefan Dziembowski and Ueli Maurer. Optimal randomizer efficiency in the bounded-storage model. *Journal of Cryptology*, 17(1):5–26, 2004.
- [DT09] Anindya De and Luca Trevisan. Extractors using hardness amplification. In *Proceedings of the 13th International Workshop on Randomization and Computation*, pages 462–475, 2009.
- [Dvi09] Zeev Dvir. Extractors for varieties. In *Proceedings of the 24th IEEE Conference on Computational Complexity*, pages 102–113, 2009.
- [GGH<sup>+</sup>07] Shafi Goldwasser, Dan Gutfreund, Alexander Healy, Tali Kaufman, and Guy Rothblum. Verifying and decoding in constant depth. In *Proceedings of the 39th ACM Symposium on Theory of Computing*, pages 440–449, 2007.
- [Gol11] Oded Goldreich. Candidate one-way functions based on expander graphs. *Studies in Complexity and Cryptography*, pages 76–87, 2011.
- [GR08a] Ariel Gabizon and Ran Raz. Deterministic extractors for affine sources over large fields. *Combinatorica*, 28(4):415–440, 2008.
- [GR08b] Venkatesan Guruswami and Atri Rudra. Explicit codes achieving list decoding capacity: Error-correction with optimal redundancy. *IEEE Transactions on Information Theory*, 54(1):135–150, 2008.
- [GRS06] Ariel Gabizon, Ran Raz, and Ronen Shaltiel. Deterministic extractors for bit-fixing sources by obtaining an independent seed. *SIAM Journal on Computing*, 36(4):1072–1094, 2006.
- [GUV09] Venkatesan Guruswami, Christopher Umans, and Salil Vadhan. Unbalanced expanders and randomness extractors from Parvaresh-Vardy codes. *Journal of the ACM*, 56(4), 2009.
- [Hås86] Johan Håstad. Almost optimal lower bounds for small depth circuits. In *Proceedings of the 18th ACM Symposium on Theory of Computing*, pages 6–20, 1986.

- [Hås87] Johan Håstad. One-way permutations in  $NC^0$ . *Information Processing Letters*, 26(3):153–155, 1987.
- [IKOS08] Yuval Ishai, Eyal Kushilevitz, Rafail Ostrovsky, and Amit Sahai. Cryptography with constant computational overhead. In *Proceedings of the 40th ACM Symposium on Theory of Computing*, pages 433–442, 2008.
- [KRVZ11] Jesse Kamp, Anup Rao, Salil Vadhan, and David Zuckerman. Deterministic extractors for small-space sources. *Journal of Computer and System Sciences*, 77(1):191–220, 2011.
- [KZ07] Jesse Kamp and David Zuckerman. Deterministic extractors for bit-fixing sources and exposure-resilient cryptography. *SIAM Journal on Computing*, 36(5):1231–1247, 2007.
- [Li11a] Xin Li. Improved constructions of three source extractors. In *Proceedings of the 26th IEEE Conference on Computational Complexity*, pages 126–136, 2011.
- [Li11b] Xin Li. A new approach to affine extractors and dispersers. In *Proceedings of the 26th IEEE Conference on Computational Complexity*, pages 137–147, 2011.
- [Lu04] Chi-Jen Lu. Encryption against storage-bounded adversaries from on-line strong extractors. *Journal of Cryptology*, 17(1):27–42, 2004.
- [LV11] Shachar Lovett and Emanuele Viola. Bounded-depth circuits cannot sample good codes. In *Proceedings of the 26th IEEE Conference on Computational Complexity*, pages 243–251, 2011.
- [MST06] Elchanan Mossel, Amir Shpilka, and Luca Trevisan. On epsilon-biased generators in  $NC^0$ . *Random Structures and Algorithms*, 29(1):56–81, 2006.
- [NZ96] Noam Nisan and David Zuckerman. Randomness is linear in space. *Journal of Computer and System Sciences*, 52(1):43–52, 1996.
- [Rao08] Anup Rao. A 2-source almost-extractor for linear entropy. In *Proceedings of the 12th International Workshop on Randomization and Computation*, pages 549–556, 2008.
- [Rao09a] Anup Rao. Extractors for a constant number of polynomially small min-entropy independent sources. *SIAM Journal on Computing*, 39(1):168–194, 2009.
- [Rao09b] Anup Rao. Extractors for low-weight affine sources. In *Proceedings of the 24th IEEE Conference on Computational Complexity*, pages 95–101, 2009.
- [Raz05] Ran Raz. Extractors with weak random seeds. In *Proceedings of the 37th ACM Symposium on Theory of Computing*, pages 11–20, 2005.
- [RRV99] Ran Raz, Omer Reingold, and Salil Vadhan. Error reduction for extractors. In *Proceedings of the 40th IEEE Symposium on Foundations of Computer Science*, pages 191–201, 1999.
- [RRV02] Ran Raz, Omer Reingold, and Salil Vadhan. Extracting all the randomness and reducing the error in Trevisan’s extractors. *Journal of Computer and System Sciences*, 65(1):97–128, 2002.

- [RY11] Ran Raz and Amir Yehudayoff. Multilinear formulas, maximal-partition discrepancy and mixed-sources extractors. *Journal of Computer and System Sciences*, 77(1):167–190, 2011.
- [RZ08] Anup Rao and David Zuckerman. Extractors for three uneven-length sources. In *Proceedings of the 12th International Workshop on Randomization and Computation*, pages 557–570, 2008.
- [Sha02] Ronen Shaltiel. Recent developments in explicit constructions of extractors. *Bulletin of the European Association for Theoretical Computer Science*, 77:67–95, 2002.
- [Sha08] Ronen Shaltiel. How to get more mileage from randomness extractors. *Random Structures and Algorithms*, 33(2):157–186, 2008.
- [Sha11] Ronen Shaltiel. An introduction to randomness extractors. In *Proceedings of the 38th International Colloquium on Automata, Languages and Programming*, pages 21–41, 2011.
- [Sho88] Victor Shoup. New algorithms for finding irreducible polynomials over finite fields. In *Proceedings of the 29th IEEE Symposium on Foundations of Computer Science*, pages 283–290, 1988.
- [SSS95] Jeanette Schmidt, Alan Siegel, and Aravind Srinivasan. Chernoff-Hoeffding bounds for applications with limited independence. *SIAM Journal on Discrete Mathematics*, 8(2):223–250, 1995.
- [Sud] Madhu Sudan. Essential coding theory — course notes. Available on the web at <http://theory.lcs.mit.edu/~madhu/coding/>.
- [TKLR09] Yael Tauman Kalai, Xin Li, and Anup Rao. 2-source extractors under computational assumptions and cryptography with defective randomness. In *Proceedings of the 50th IEEE Symposium on Foundations of Computer Science*, pages 617–626, 2009.
- [Tre01] Luca Trevisan. Extractors and pseudorandom generators. *Journal of the ACM*, 48(4):860–879, 2001.
- [TV00] Luca Trevisan and Salil Vadhan. Extracting randomness from samplable distributions. In *Proceedings of the 41st IEEE Symposium on Foundations of Computer Science*, pages 32–42, 2000.
- [Vad] Salil Vadhan. Pseudorandomness. *Foundations and Trends in Theoretical Computer Science (to appear)*.
- [Vad04] Salil Vadhan. Constructing locally computable extractors and cryptosystems in the bounded-storage model. *Journal of Cryptology*, 17(1):43–77, 2004.
- [Vio10] Emanuele Viola. The complexity of distributions. In *Proceedings of the 51st IEEE Symposium on Foundations of Computer Science*, pages 202–211, 2010.
- [Vio11] Emanuele Viola. Extractors for circuit sources. In *Proceedings of the 52nd IEEE Symposium on Foundations of Computer Science (to appear)*, 2011.

- [Yao85] Andrew Yao. Separating the polynomial-time hierarchy by oracles. In *Proceedings of the 26th IEEE Symposium on Foundations of Computer Science*, pages 1–10, 1985.
- [Yeh11] Amir Yehudayoff. Affine extractors over prime fields. *Combinatorica*, 31(2):245–256, 2011.
- [Zim10] Marius Zimand. Simple extractors via constructions of cryptographic pseudo-random generators. *Theoretical Computer Science*, 411(10):1236–1250, 2010.