

Small-bias probability spaces: efficient constructions and applications *

Joseph Naor †

Computer Science Department
Technion
Haifa 32000, Israel

Moni Naor

IBM Almaden Research Center
650 Harry Road
San Jose, CA 95120, USA

Abstract

We show how to efficiently construct a small probability space on n binary random variables such that for every subset, its parity is either zero or one with “almost” equal probability. They are called ϵ -biased random variables. The number of random bits needed to generate the random variables is $O(\log n + \log \frac{1}{\epsilon})$. Thus, if ϵ is polynomially small, then the size of the sample space is also polynomial. Random variables that are ϵ -biased can be used to construct “almost” k -wise independent random variables where ϵ is a function of k .

These probability spaces have various applications:

1. Derandomization of algorithms: many randomized algorithms that require only k -wise independence of their random bits (where k is bounded by $O(\log n)$), can be derandomized by using ϵ -biased random variables.
2. Reducing the number of random bits required by certain randomized algorithms, e.g., verification of matrix multiplication.
3. Exhaustive testing of combinatorial circuits. We provide the smallest known family for such testing.
4. Communication complexity: two parties can verify equality of strings with high probability exchanging only a logarithmic number of bits.
5. Hash functions: we can construct a polynomial sized family of hash functions such that with high probability, the sum of a random function over two different sets is not equal.

*A preliminary version of this paper appeared in the Proceedings of the 22nd Annual ACM Symposium on Theory of Computing, Baltimore, Maryland (1990), pp. 213-223.

†Most of this work was done while the author was at the Computer Science Department, Stanford University, Stanford, CA and supported by contract ONR N00014-88-K-0166.

1 Introduction

Randomness plays a significant role in computer science. However, it is often desirable to reduce the amount of randomness required. The purpose of this paper is to construct small probability spaces that approximate larger ones. Let x_1, \dots, x_n be $\{0, 1\}$ Bernoulli random variables and let Ω be the probability space associated with them. If the random variables are independent, then Ω contains all 2^n possible assignments. Our goal is to construct a much smaller probability space that will behave similarly to Ω in certain respects. Such small probability spaces have proved to be very useful.

One of the main approaches taken by previous researchers to reduce the size of the sample space was to require only limited independence among the random variables (as opposed to full independence). Our approach is different; it is based on the equivalence of the following two conditions: (See [18] and [52]).

1. The random variables are independent and for all i , $\text{Prob}[x_i = 0] = \text{Prob}[x_i = 1]$.
2. For every subset $S \subseteq \{1, \dots, n\}$, it is equally likely that the parity of the subset (i.e., the number of “ones”) is either zero or one.

We are going to relax the second condition and construct a probability distribution such that for every subset $S \subseteq \{1, \dots, n\}$, it is “almost” equiprobable that the parity of the subset is zero or one. To be more precise, we require that for every subset S ,

$$\left| \text{Prob}\left[\sum_{i \in S} x_i = 0\right] - \text{Prob}\left[\sum_{i \in S} x_i = 1\right] \right| \leq \epsilon$$

This quantity was called by Vazirani [52] the *bias* of the subset S . The cardinality of the sample space we construct is $2^{O(\log \frac{1}{\epsilon} + \log n)}$. Hence, if ϵ is typically polynomially small, then the size of the sample space is also polynomial.

We also define random variables that are k -wise ϵ -biased. For them, only the bias of subsets that are smaller than k is guaranteed to be bounded by ϵ . We present a more efficient construction for such random variables: the logarithm of the cardinality of the sample space is $O(\log k + \log \log n + \log \frac{1}{\epsilon})$. In Section 5 we show how the k^{th} moment of the sum of k -wise ϵ -biased random variables can be bounded via the k^{th} moment of the binomial distribution on uniform independent Bernoulli variables. This is an important tool for analyzing the behavior of the distribution of the sum.

We can use k -wise ϵ -biased random variables to construct k -wise δ -dependent random variables, i.e., variables such that the variation distance between the distribution of any subset of

k variables and the uniform distribution (on k variables) is at most δ . Their construction is described in Section 4.

The ϵ -biased random variables are constructed in three stages. In the first stage, we construct a sample space \mathcal{F} of random variables such that the bias of every subset S is bounded by some constant. Sampling from \mathcal{F} requires $O(\log n)$ random bits. In the second stage, \mathcal{F} is sampled l times (not necessarily independently), where l depends on ϵ . In the third stage, the ϵ -biased random variables are generated by picking a linear combination of the assignments sampled in the second stage. Constructing the probability space is described in Section 3.

Another interpretation of our result is via Fourier transforms. (See Section 2 for precise definitions). For the uniform distribution, all the coefficients of its Fourier transform are zero, except for the free coefficient. For an ϵ -biased distribution, the absolute value of each coefficient is at most $\frac{\epsilon}{2^n}$.

Derandomizing algorithms has attracted much attention in recent years. For the purpose of derandomization, our distribution can replace the uniform one in many cases, so as to allow an exhaustive search for a good point in a polynomial-size sample space. We exemplify this by providing a *polynomial* size sample space for the set balancing problem; this also yields an NC¹ algorithm. A previous approach to derandomizing the set balancing problem [13, 39] was to first construct an $n^{O(\log n)}$ sample space, and then conduct a binary search for a good point. As a result, their time bounds are worse. Another problem we address is finding a heavy codeword in a linear code. Using ϵ -biased random variables, we provide the first NC algorithm to the problem.

Another application of our probability distribution is for reducing the number of random bits required for several randomized algorithms. Karp and Pippenger [32] suggested that random bits can be viewed as a resource (just as time and space) which is best to use as little as possible. One motivation to consider randomness as a resource is practical. Random bits are hard to produce and devices that generate them, such as Geiger counters and Zener diodes are slow. Another reason is from the complexity theoretic point of view: to provide a full scale of options between an algorithm that is completely deterministic, and a randomized algorithm that consumes many bits.

Examples of previous work in reducing the number of random bits in randomized algorithms are [2, 11, 20, 28, 32, 33, 43, 47, 49]. Karloff and Raghavan [33] for example, studied several randomized algorithms for selection and sorting and showed that they can be run successfully when only $O(\log n)$ random bits are available.

In Section 7 we describe how to reduce the number of random bits for three problems. The first one is matrix multiplication verification. Given three $n \times n$ matrices A , B and C , how

can one verify that $A \cdot B = C$ without resorting to matrix multiplication. We show how to do that in $O(n^2)$ time using $O(\log n)$ random bits, thus improving on a previous algorithm of [26] that required $O(n)$ random bits. We next show how Adi Shamir's Boolean matrix multiplication algorithm can be implemented by using $O(\log n)$ random bits without an increase in the probability of error. The third problem is verifying n equalities of the form $a^{x_i} = y_i \bmod p$ where p is a prime. We show how to do that using only $O(\log n)$ random bits and n multiplications, instead of $n \log p$ multiplications needed for testing each equality separately.

In Section 8 we show how to apply our construction to generate fault-diagnostic tests for combinatorial circuits. A well known problem in that area is to construct a small collection of assignments to inputs of a circuit such that for any k inputs, all possible configurations appear. Using k -wise δ -dependent probability spaces, we can get the best explicit constructions. These are optimal up to the constant factor in the exponent.

Our techniques can be used to minimize the communication complexity of protocols for testing equality of two strings, while achieving a very low probability of error. Similarly, the techniques can be applied to construct a small family of hash functions with the property that summing the hash function over different sets yields a different value with high probability. These two applications are discussed in Section 9.

In Section 10 we briefly survey recent papers that have used the constructions described in this paper since it appeared in [42].

Independent of our work, Peralta [44] has considered ϵ -bias probability spaces as well, and showed some applications to number theoretic algorithms. His construction is based on quadratic residues and Weil's Theorem.

2 Preliminaries and definitions

Let $\mathbf{x} = x_1, \dots, x_n$ be $\{-1, 1\}$ random variables and D their joint probability distribution.

Definition 2.1 *The bias of a subset $S \subseteq \{1, \dots, n\}$ for a distribution D is defined to be,*

$$\text{bias}_D(S) = \left| \text{Prob}_D \left[\prod_{i \in S} x_i = -1 \right] - \text{Prob}_D \left[\prod_{i \in S} x_i = 1 \right] \right|$$

Definition 2.2 *The variables x_1, \dots, x_n are ϵ -biased if for all S , $\text{bias}_D(S) \leq \epsilon$. They are said to be k -wise ϵ -biased if for all subsets S such that $|S| \leq k$, $\text{bias}_D(S) \leq \epsilon$.*

Let U denote the uniform distribution and $D(S)$ the distribution D restricted to the subset S . The variation distance between two distributions D_1 and D_2 defined over the same probability

space Ω is

$$\|D_1 - D_2\| = \sum_{\omega \in \Omega} |D_1(\omega) - D_2(\omega)|$$

Definition 2.3 *The variables x_1, x_2, \dots, x_n are defined to be k -wise δ -dependent if for all subsets S such that $|S| \leq k$,*

$$\|U(S) - D(S)\| \leq \delta$$

A similar definition was made by R. Ben-Nathan [12].

The set $f : \{-1, 1\}^n \rightarrow \mathfrak{R}$ of real functions on the n -dimensional cube forms a 2^n -dimensional real vector space. The inner product of two functions f and g is defined as,

$$\langle f, g \rangle = \frac{1}{2^n} \sum_{x \in \{-1, 1\}^n} f(x)g(x)$$

The basis of this vector space is given by the following family of functions, the characters of Z_2^n . Define for each subset $S \subseteq \{1 \dots n\}$,

$$\chi_S(x_1, \dots, x_n) = \prod_{i \in S} x_i$$

The above basis has the following properties:

1. For all $S \subseteq \{1 \dots n\}$, the family $\{\chi_S\}$ forms an orthonormal basis: if $S_1 \neq S_2$, then $\langle \chi_{S_1}, \chi_{S_2} \rangle = 0$, and for every S , $\langle \chi_S, \chi_S \rangle = 1$.
2. For every S_1, S_2 : $\chi_{S_1} \cdot \chi_{S_2} = \chi_{S_1 \Delta S_2}$, where $S_1 \Delta S_2$ is the symmetric difference of S_1 and S_2 .

The Fourier transform of a function f is its expansion as a linear combination of the χ_S 's. Every function has a unique expression and the coefficients in this expression are called the Fourier coefficients, denoted by $\hat{f}(S)$ (for $S \subseteq \{1 \dots n\}$). Hence, $f = \sum_S \hat{f}(S) \chi_S$.

Since the family χ_S forms an orthonormal basis, Fourier coefficients are found via:

$$\hat{f}(S) = \langle f, \chi_S \rangle$$

For a probability distribution D , $D = \sum_S \hat{D}(S) \chi_S$. The following theorem was proved by Diaconis and Shahshahani (see [22, Lemma 1, p. 24]). We provide a proof for the sake of completeness. (A slightly weaker bound was actually proved by Vazirani [52]; it motivated us to consider ϵ -biased probability spaces).

Theorem 2.1 Let x_1, \dots, x_n be $\{-1, 1\}$ random variables and let D be their joint probability distribution. Then,

$$\|D - U\| \leq 2^n \cdot \left(\sum_{S \subseteq \{1 \dots n\}} \hat{D}(S)^2 \right)^{\frac{1}{2}} = \left(\sum_{S \subseteq \{1 \dots n\}} \text{bias}_D^2(S) \right)^{\frac{1}{2}}$$

Proof: We first evaluate $\hat{D}(\emptyset)$,

$$\hat{D}(\emptyset) = \langle D, \chi_\emptyset \rangle = \sum_{x \in \{-1, 1\}^n} \frac{D(x)}{2^n} = \frac{1}{2^n}.$$

The square of the variation distance is,

$$\|D - U\|^2 = \left(\sum_{x \in \{-1, 1\}^n} \left| D(x) - \frac{1}{2^n} \right| \right)^2$$

By the Cauchy-Schwarz inequality,

$$\left(\sum_{x \in \{-1, 1\}^n} \left| D(x) - \frac{1}{2^n} \right| \right)^2 \leq 2^n \cdot \sum_{x \in \{-1, 1\}^n} \left(D(x) - \frac{1}{2^n} \right)^2 = 2^n \cdot \sum_{x \in \{-1, 1\}^n} \left(\sum_S \hat{D}(S) \cdot \chi_S(x) - \frac{1}{2^n} \right)^2$$

Since $\hat{D}(\emptyset) = \frac{1}{2^n}$,

$$\begin{aligned} 2^n \cdot \sum_{x \in \{-1, 1\}^n} \left(\sum_S \hat{D}(S) \cdot \chi_S(x) - \frac{1}{2^n} \right)^2 &= 2^n \cdot \sum_{x \in \{-1, 1\}^n} \left(\sum_{S \neq \emptyset} \hat{D}(S) \cdot \chi_S(x) \right)^2 = \\ 2^n \cdot \sum_{x \in \{-1, 1\}^n} \sum_{S_1 \neq \emptyset} \sum_{S_2 \neq \emptyset} \hat{D}(S_1) \hat{D}(S_2) \chi_{S_1}(x) \chi_{S_2}(x) &= 2^n \cdot \sum_{S_1 \neq \emptyset} \sum_{S_2 \neq \emptyset} \hat{D}(S_1) \hat{D}(S_2) \sum_{x \in \{-1, 1\}^n} \chi_{S_1}(x) \chi_{S_2}(x) = \\ &= 2^{2n} \cdot \sum_{S_1 \neq \emptyset} \sum_{S_2 \neq \emptyset} \hat{D}(S_1) \hat{D}(S_2) \langle \chi_{S_1}, \chi_{S_2} \rangle \end{aligned}$$

If $S_1 \neq S_2$, then $\langle \chi_{S_1}, \chi_{S_2} \rangle = 0$; otherwise, $\langle \chi_{S_1}, \chi_{S_2} \rangle = 1$. Hence,

$$2^{2n} \cdot \sum_{S_1 \neq \emptyset} \sum_{S_2 \neq \emptyset} \hat{D}(S_1) \hat{D}(S_2) \langle \chi_{S_1}, \chi_{S_2} \rangle = 2^{2n} \cdot \sum_S \hat{D}(S)^2$$

Which implies that,

$$\|D - U\| \leq 2^n \cdot \left(\sum_{S \subseteq \{1 \dots n\}} \hat{D}(S)^2 \right)^{\frac{1}{2}}$$

To complete the proof, recall that the definition of the bias of a subset S is,

$$\begin{aligned} \text{bias}_D(S) &= |\text{Prob}_D[\chi_S(x) = -1] - \text{Prob}_D[\chi_S(x) = 1]| = \\ &= \left| \sum_{x \in \{-1, 1\}^n} D(x) \chi_S(x) \right| = 2^n | \langle D, \chi_S \rangle | = |2^n \hat{D}(S)| \end{aligned}$$

and hence,

$$2^n \cdot \left(\sum_{S \subseteq \{1 \dots n\}} \hat{D}(S)^2 \right)^{\frac{1}{2}} = \left(\sum_{S \subseteq \{1 \dots n\}} \text{bias}_D^2(S) \right)^{\frac{1}{2}}$$

□

Corollary 2.1 *If the random variables x_1, \dots, x_n are ϵ -biased with respect to a distribution D , then they are also k -wise δ -dependent, for $\delta = 2^{\frac{k}{2}} \cdot \epsilon$.*

A *binary linear code* C is a linear subspace of $\{0, 1\}^m$. If C has dimension n then C is called an $[m, n]$ code and it transforms words of length n into codewords of length m . A *generator matrix* G for a linear code C is an $n \times m$ matrix for which the rows are a basis of C . If G is a generator matrix for C , then the code can be defined as

$$C = \{aG \mid a \in \{0, 1\}^n\}$$

The distance between two codewords is defined to be their *Hamming distance*. $C[m, n, \beta]$ denotes an $[m, n]$ code C for which the distance between any two codewords is at least $\beta \cdot m$. The *weight* of a codeword is the number of non-zero symbols that it contains. For more details on linear codes the reader is referred to [40].

3 Constructing the probability distribution

In this section we show how to construct a probability space Ω of $\{0, 1\}$ random variables $\mathbf{x} = x_1, \dots, x_n$ which are ϵ -biased. The cardinality of the sample space will be $2^{O(\log \frac{1}{\epsilon} + \log n)}$. The joint probability distribution of the variables is denoted by D . We define the function $\chi'_S(x_1, \dots, x_n)$ for $\{0, 1\}$ random variables as follows. For each subset $S \subseteq \{1 \dots n\}$,

$$\chi'_S(x_1, \dots, x_n) = \sum_{i \in S} x_i \pmod{2}.$$

The construction consists of three stages:

1. A polynomial size family \mathcal{F} of $\{0, 1\}^n$ vectors is generated with the following property. Let r be a vector chosen from \mathcal{F} uniformly at random. For all subsets $S \subseteq \{1, \dots, n\}$,

$$\text{Prob}[\chi'_S(r) = 1] \geq \beta$$

where β is some constant. Constructing such a family is discussed in Section 3.1.

2. The vectors r_1, \dots, r_l are sampled from \mathcal{F} (not necessarily independently, but via a Markov process) such that for all subsets $S \subseteq \{1, \dots, n\}$,

$$\text{Prob}[\text{For all } i, 1 \leq i \leq l, \chi'_S(r_i) = 0] \leq \epsilon$$

The value of l will turn out to be $O(\log \frac{1}{\epsilon})$. Sampling the family \mathcal{F} is discussed in Section 3.2.

3. The assignment to the random variables x_1, \dots, x_n is a combination of the vectors sampled at the previous stage. Let $\vec{a} = (a_1, \dots, a_l)$ be chosen uniformly at random from $\{0, 1\}^l$. Then,

$$\mathbf{x} = \sum_{i=1}^l a_i r_i$$

In Section 3.3 we discuss how to choose \vec{a} so that \mathbf{x} is ϵ -biased.

3.1 Constructing the family \mathcal{F}

We will need a family \mathcal{F} with the above properties for constructing ϵ -biased probability spaces and also for other applications as well. (See Section 7). For the latter applications, we extend the requirements from \mathcal{F} to any ring: given a vector v of elements in the ring, a vector r such that $\langle v \cdot r \rangle \neq 0$ is called a *distinguisher* with respect to v . The goal is to find a small collection of vectors (the family \mathcal{F}) such that for any non-zero vector v , if r is chosen at random, then $\text{Prob}[\langle v \cdot r \rangle \neq 0] \geq \beta$. If the ring is $\text{GF}[2]$, this requirement is exactly the one mentioned above, where v is the characteristic vector of the subset S . Henceforth, we describe the construction for $\text{GF}[2]$. It can easily be generalized for any ring by substituting 1 by some non-zero element of the ring.

We present two methods for constructing the family \mathcal{F} . The first method (Section 3.1.1) can be applied to any ring, whereas the second one (Section 3.1.2) is applicable only to $\text{GF}[2]$. Another advantage of the first method is that computing the value of a random variable $x_i \in \mathbf{x}$ can be done in $O(1)$ operations on words of length $\log n$. On the other hand, the second method provides a general context, i.e., *linear codes*.

Proposition 3.1 *Suppose that r is a vector chosen uniformly at random from $\{0, 1\}^n$. Then for all vectors $v \in \{0, 1\}^n$, $v \neq \vec{0}$, $\text{Prob}[\langle v \cdot r \rangle \neq 0] \geq \frac{1}{2}$.*

Proof: Let j denote the number of non-zero entries in v . The number of different vectors r such that $v \cdot r = 0$ is at most $2^{j-1} \cdot 2^{n-j} = 2^{n-1}$, whereas the number of distinct choices for r is 2^n . Hence, the probability that $v \cdot r = 0$ is at most $\frac{1}{2}$. \square

Unfortunately, this method of generating a distinguisher requires n random bits. Hence, our aim is to show that a distinguisher can be generated with probability at least β using only $O(\log n)$ random bits. This will guarantee that the size of \mathcal{F} is polynomial. Note that \mathcal{F} must contain at least n vectors (potential distinguishers). Otherwise, the rank of the matrix whose columns are the distinguishers is less than n .

For our purposes, a collection \mathcal{F} of vectors has to be constructed such that:

1. The size of \mathcal{F} is “small”.
2. It is “easy” to sample uniformly from \mathcal{F} .
3. Given any non-zero vector v , a constant fraction of the vectors in \mathcal{F} are distinguishers with respect to v .

3.1.1 Constructing a small set of distinguishers

A natural approach to the problem of reducing the number of random bits in a randomized algorithm is to show that limited independence of the random variables suffices to assure a high probability of success. (See e.g., [5, 36, 38]). However, in our case, it is not clear from the proof of Proposition 3.1 how many vectors r remain distinguishers with respect to the vector v when the entries of R are not completely independent. Moreover, an example can be constructed in which if the entries of r are chosen pairwise independently, no distinguisher will be generated.

Though limited independence is not sufficient for our purposes, we make use of it in two ways: one is that suggested in [17] that if we sample a universe looking for elements of some fixed subset, and if the expected number of elements we hit is greater than 1, then by making our choices only pairwise independent, we are not decreasing the chances of hitting an element of the subset by much. The other is that if the vector v has at most c non-zero entries, then the elements of r can be chosen c -wise independent, and with probability at least $\frac{1}{2}$ r is a distinguisher.

We now describe how the above-mentioned difficulties for generating distinguishers can be overcome. In what follows we will need n random variables such that:

1. Each random variable is uniformly distributed in $\{1, \dots, n\}$.
2. Every subset of the random variables of cardinality at most c is independent.

There are known methods of generating such random variables that use only $O(c \log n)$ random bits. ([38], [5], [17]).

We first assume that l , the precise number of non-zero elements in v , is known to be in the range $[k \dots 2k - 1]$. A two-step process is applied.

1. The vector v is replaced by a new vector $v' = (v'_1, \dots, v'_n)$ that contains only c (for some constant) non-zero elements. (Any non-zero element of v' is also non-zero in v).

2. Now, the elements of the vector r can be chosen c -wise independent, yet Proposition 3.1 still holds.

As we do not have direct access to the elements of v , we show instead how to emulate the above Step 1 with high probability. Let $u = (u_1, \dots, u_n)$ and $w = (w_1, \dots, w_n)$ be two random vectors such that:

1. The entries of u are c -wise independent (c is a constant whose value will be specified later) where for all $1 \leq i \leq n$, $\text{Prob}[u_i = 0] = \text{Prob}[u_i = 1] = \frac{1}{2}$.
2. The entries of w are pairwise independent where for all $1 \leq i \leq n$, $\text{Prob}[w_i = 1] = \frac{2}{k}$. (If $k = 1$, then $\text{Prob}[w_i = 1] = 1$).

We can assume w.l.o.g that $k|n$. To generate the vector w we generate n random variables z_1, z_2, \dots, z_n that are pairwise independent, and each is uniformly distributed in $\{1, \dots, n\}$. We then set w_i to 1 if $z_i \leq \frac{2n}{k}$.

Let us now define the random vector $r = (r_1, \dots, r_n)$ that will be used as a distinguisher. For all $1 \leq i \leq n$,

$$r_i = \begin{cases} 1 & \text{if } u_i = 1 \text{ and } w_i = 1 \\ 0 & \text{otherwise} \end{cases}$$

Lemma 3.1 *The above vector r is a distinguisher with probability at least $\frac{1}{4}$ for any vector v for which l , the number of non-zero elements, is in the range $[k, 2k]$.*

Proof: Let us define the vector v' (from Step 1): for all $1 \leq i \leq n$,

$$v'_i = \begin{cases} 1 & \text{if } v_i = 1 \text{ and } w_i = 1 \\ 0 & \text{otherwise} \end{cases}$$

It is clear that the lemma will follow if we show that the vector u is a distinguisher with respect to v' with probability at least $\frac{1}{4}$. To do that, it suffices to prove that v' will contain at least one non-zero element of v , and at most c non-zero elements of v with probability at least $\frac{1}{2}$. As the elements of the vector u are c -wise independent, the proof of Proposition 3.1 still holds when the number of non-zero elements is less than c .

Generating the vector v' can be thought of as a binomial random variable where each non-zero entry of v decides with probability p (to be specified later) whether it remains non-zero in v' . The random choices are pairwise independent. Let h be a random variable that denotes the number of non-zero elements in v' . It is well known that $E[h] = pl$ and $\text{Var}[h] = p(1-p)l$. We chose the value of p such that $pk = 2$.

Claim: $\text{Prob}[0 < h \leq 7] \geq \frac{1}{2}$

We prove the claim by Chebyshev's inequality [Fe] which states that

$$\text{Prob}[|X - E[X]| \geq \lambda] \leq \frac{\text{Var}[X]}{\lambda^2}$$

where X is a random variable. It is enough to verify the claim in the two extreme cases when $l = k$ and $l = 2k$. Thus, substituting $\lambda = 2$ and $\lambda = 3$, we get that

$$\text{Prob}[|h - pk| \geq 2] \leq \frac{pk(1-p)}{4} \leq \frac{1}{2}$$

$$\text{Prob}[|h - 2pk| \geq 3] \leq \frac{2pk(1-p)}{9} \leq \frac{4}{9}$$

Hence, we can choose $c = 2kp + 3 = 7$ and this is enough to insure success with probability at least $\frac{1}{2}$. \square .

We conclude that if the approximate number of non-zero entries in v is known, then $O(\log n)$ random bits suffice to insure high probability of success. What can we do if this is not known? We follow the above algorithm and construct $\log n$ collections $\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_{\log n}$, where \mathcal{F}_i is generated under the assumption that the number of non-zero entries in v is between 2^{i-1} and 2^i . Lemma 3.1 implies that at least $\frac{1}{4}$ of the members of at least one collection will be distinguishers.

The same random bits can be used to sample the $\log n$ collections, and we obtain a set of $\log n$ vectors $r^1, r^2, \dots, r^{\log n}$ such that at least one of them is a distinguisher with probability at least $\frac{1}{4}$. Instead of testing each vector separately, we can generate from them a single vector that is a distinguisher with probability at least $\frac{1}{8}$.

Let S be a subset of $\{1, \dots, \log n\}$ which is chosen uniformly at random and let r be defined by

$$r = \sum_{i \in S} r^i$$

Lemma 3.2 *The vector r is a distinguisher with respect to v with probability at least $\frac{1}{8}$.*

Proof: We need the following claim:

Claim: Let v be a vector such that the vector r_1 is a distinguisher and the vector r_2 is not a distinguisher with respect to v . Then the vector $r_1 + r_2$ is a distinguisher with respect to v .

Proof: $(r_1 + r_2) \cdot v = r_1 \cdot v + r_2 \cdot v = r_1 \cdot v \neq 0$.

Let the number of non-zero entries in v be between 2^{j-1} and 2^j and assume that r^j is a distinguisher. This will happen with probability at least $\frac{1}{4}$. If $\sum_{i \in S - \{j\}} r^i$ is not a distinguisher, then with probability $\frac{1}{2}$, $j \in S$ and according to the above claim, r will be a distinguisher.

Otherwise, again with probability $\frac{1}{2}$, $j \notin S$, and according to the above claim, r will be a distinguisher. \square

To summarize, we describe the algorithm to generate a random vector r . We assume that n is a power of 2.

1. Generate the following random variables:
 - (a) z_1, z_2, \dots, z_n : n random variables that are pairwise independent and uniformly distributed in $\{1, \dots, n\}$.
 - (b) $u = (u_1, \dots, u_n)$ a vector whose entries are 7-wise independent and uniformly distributed in $\{0, 1\}$.
 - (c) A random subset $S \subseteq \{1, \dots, \log n\}$.
2. For each $1 \leq i \leq n$, compute $j_i = \log n - \max\{j | 2^j < z_i\}$.
3. For $1 \leq l \leq \log n$, compute $c_l = |S \cap \{1, \dots, l\}| \cdot 1$. (This is scalar multiplication in the ring, e.g., in $\text{GF}[2]$ it is $|S \cap \{1, \dots, l\}| \bmod 2$).
4. For each $1 \leq i \leq n$, compute:

$$r_i = \begin{cases} c_{j_i} & \text{if } u_i = 1 \\ 0 & \text{otherwise} \end{cases}$$

Theorem 3.1 *The algorithm described above uses $O(\log n)$ random bits and generates a distinguisher with probability at least $\beta = \frac{1}{8}$. For all i , $1 \leq i \leq n$, the complexity of computing the value of the random variable x_i is $O(1)$ operations on words of size $O(\log n)$.*

Proof: Steps 1a, 1b and 1c require each $O(\log n)$ bits. (In fact, Steps 1b and 1c can be implemented recursively using $O(\log \log n)$ bits). Given i , to compute r_i , one needs to know the value of z_i, c_{j_i} and u_i ; that requires a constant number of operations. In Step 3, computing c_{j_i} requires counting the number of 1-s in a word of length $\log n$ that describes S . \square

3.1.2 A construction based on linear codes

Here we describe how to generate a family \mathcal{F} of size $O(n)$ via linear codes due to J. Bruck (private communication). The construction works for $\text{GF}[2]$. Let C be a linear code. The weight of a codeword is defined to be the number of non-zero entries.

Proposition 3.2 *The minimum distance of a linear code is equal to the minimum weight of a codeword.*

Suppose we have a linear code $C[n, m, \beta]$, i.e. it maps $\{0, 1\}^n$ words into $\{0, 1\}^m$ codewords and its minimum distance is $\beta \cdot m$. Linear codes for which $m = O(n)$ and β is some constant exist and are constructible. (For example, Justesen codes [30], [40]). Let G be the generator matrix of this code. For any non-zero $\{0, 1\}^n$ vector v , $v \cdot G$ contains at least βm non-zero entries. (By the above proposition). Hence, if we choose a column in G uniformly at random, $\text{Prob}[\langle v \cdot r \rangle = 1] \geq \beta$. The family \mathcal{F} is the set of columns in G . The relation with linear codes holds in the other direction as well. Given a family \mathcal{F} , consider its members as columns of a generator matrix of a linear code. It follows from the proposition that the minimum distance of this code is $\beta|\mathcal{F}|$. Hence, the construction in Section 3.1.1 can be regarded as a linear code.

Given that good codes exist, what advantages does the first method have? The first method has the property that it works for more general cases where the function χ'_S is defined on any group, not just addition modulo 2. This will be used in Section 7.1. Another advantage is in computing a single entry of the sampled vector r . This is very simple in the first method (Theorem 3.1) whereas all known methods for using (traditional linear codes) are more complicated and require exponentiation.

Using our techniques we can actually enhance the error-correction of a linear code without decreasing the rate by much. This is further investigated in [6].

3.2 Sampling the family \mathcal{F}

The problem of obtaining the vectors r_1, \dots, r_l in Stage 2 with the desired property can be abstracted in the following way. Suppose that there is a universe, and we wish to find a member in a certain subset \mathcal{S} of it. (In our case it is the set of vectors for which $\chi_S = 1$). Suppose also that we have a sampling algorithm that uses k random bits and has probability β of picking a member of the desired set. By sampling l times independently, a set of l elements is generated, such that with probability greater than $1 - \beta^{-l}$, at least one of them is a member of the desired set. A straightforward implementation would require kl random bits.

Several papers have addressed the question of achieving the probability error while requiring fewer random bits [1, 20, 28, 32, 47, 49]. We consider the method of [1] which is used by [20, 28]: For a graph $G = (V, E)$, consider any 1-1 correspondence $f : V \rightarrow \{0, 1\}^k$, the different assignments to the random bits of the sampling algorithm. To generate the l samples, choose a random vertex of G and perform a random walk of length l . Each vertex in the random walk corresponds to a sample point. The number of random bits required is the sum of those needed for sampling one vertex, and those needed for performing the random walk.

Let λ_0 be the largest eigenvalue of G and $\bar{\lambda}$ be the eigenvalue of second largest value in G . Let α denote the percentage of vertices that are not members in \mathcal{S} . Cohen and Wigderson [21,

Theorem 4.5] show that if G is a d -regular graph such that

$$2\alpha \leq \left(\frac{\bar{\lambda}}{\lambda_0}\right)^2$$

then the probability that at least one of the l samples is a member of \mathcal{S} is at least

$$1 - \left(\frac{\sqrt{2}|\bar{\lambda}|}{\lambda_0}\right)^l.$$

Constructions for regular graphs of constant degree such that the second eigenvalue is bounded away from λ_0 are known [27, 29, 37]. For degree regular graphs of degree d , $\lambda_0 = d$ and the value of $\bar{\lambda}$ can be almost $2\sqrt{d}$. For a given expander G , let

$$\alpha_G = \frac{1}{2} \left(\frac{\bar{\lambda}}{\lambda_0}\right)^2$$

and let

$$l = \frac{\log \frac{1}{\epsilon}}{\log \left(\frac{\sqrt{2}|\bar{\lambda}|}{\lambda_0}\right)}.$$

In our case, β , the percentage of good vectors, is too small to apply the method directly and we need some initial amplification. To do that, each vertex would now correspond to a set of assignments to the random bits, such that the probability that at least one element associated with a randomly chosen vertex is a member of \mathcal{S} , is at least $1 - \alpha_G$. For the purposes of this paper, this can be done by letting each vertex correspond to h independent samples where $h = \log_{1-\beta} \alpha_G$.

To conclude, the number of random bits used is $\log |\mathcal{F}| + l \cdot \log d$.

3.3 Combining the samples

We have to specify how to choose a linear combination \vec{a} of the vectors r_1, \dots, r_l , given that for any subset S , $\text{Prob}[\text{for all } i, \chi'_S(r_i) = 0] \leq \epsilon$. The simplest way to select \vec{a} is to choose it uniformly at random.

Claim 3.1 *The random variables x_1, \dots, x_n generated by choosing \vec{a} uniformly at random from $\{0, 1\}^l$ are ϵ -biased.*

Proof: Let S be any subset of $\{1 \dots n\}$. If there exists a vector r_j among the vectors r_1, \dots, r_k such that $\chi'_S(r_j) = 1$, then by arguments similar to those of Lemma 3.2,

$$\text{Prob}_D[\chi'_S(x_1, \dots, x_n) = 1] = \text{Prob}_D[\chi'_S(x_1, \dots, x_n) = 0]$$

Since the probability of this happening is at least $1 - \epsilon$, $\text{bias}_D(S) \leq \epsilon$. \square

The number of random bits required is l . To conclude,

Theorem 3.2 *Generating n $\{0, 1\}$ random variables that are ϵ -biased can be done using $O(\log n + \log \frac{1}{\epsilon})$ random bits. Thus, the size of the sample space is $2^{O(\log n + \log \frac{1}{\epsilon})}$. Given the random bits, computing the value of a random variable can be done in time polylogarithmic in n .*

4 Generating k -wise δ -dependent random variables

We are insured by Corollary 2.1 that if the random variables x_1, \dots, x_n are ϵ -biased, then they are k -wise δ -dependent for $\delta = 2^{\frac{k}{2}} \cdot \epsilon$. However, there is a more efficient construction. This can be done by combining our methods with those of [5] for generating k -wise independent variables.

Suppose we want to generate $\{0, 1\}$ uniform random variables y_1, \dots, y_n that are k -wise independent. [5] suggest that this can be done by taking n vectors L_1, \dots, L_n of length h such that the vectors are k -wise linearly independent over $\text{GF}[2]$. If the vectors L_1, \dots, L_n are columns of the parity check matrix of BCH codes, then h is $\frac{k}{2} \log n$. Let R be a vector chosen uniformly at random from $\{0, 1\}^h$; for all i , $1 \leq i \leq n$, let $y_i = L_i \cdot R$. The number of random bits required for the construction is $k \log n$.

In order to improve on Corollary 2.1, instead of choosing R uniformly at random, suppose that the entries of R are ϵ -biased random variables.

Lemma 4.1 *Let y_1, \dots, y_n be random variables generated by the above method, where the entries of R are ϵ -biased random variables. Then, y_1, \dots, y_n are k -wise ϵ -biased.*

Proof: Let $S \subseteq \{1, \dots, n\}$ be a subset of cardinality at most k . We bound $\text{bias}(S)$. For all $i \in S$, $y_i = L_i \cdot R$. Hence,

$$\sum_{i \in S} y_i = \sum_{i \in S} L_i \cdot R = R \cdot \sum_{i \in S} L_i = R \cdot M$$

For $i \in S$, the vectors L_i are linearly independent, and hence $M \neq 0$ and $\text{bias}(S) \leq \epsilon$. \square

The improvement over Corollary 2.1 in the cardinality of the sample space is that now we have decreased the number of ϵ -biased random variables from n to $k \log n$. Recall that random variables that are k -wise ϵ -biased, are also k -wise $2^{\frac{k}{2}} \cdot \epsilon$ -dependent.

Lemma 4.2 *The logarithm of the cardinality of the sample space needed for constructing k -wise δ -dependent random variables is $O(k + \log \log n + \log \frac{1}{\delta})$.*

5 A moment inequality

Basic tools in probabilistic analysis are moment inequalities [25] that bound the deviation of a random variable from its expected value. More specifically, let y be a random variable such that

$E[y] = 0$, then,

$$\text{Prob}[|y| \geq \lambda] \leq \frac{E[|y|^k]}{\lambda^k}$$

Let b_1, \dots, b_n be random variables that have a binomial distribution, i.e., they are independent and take their values from $\{-1, 1\}$ uniformly. The k^{th} moment of their sum will be denoted by B_k , that is $B_k = E[|b_1 + \dots + b_n|^k]$.

We formulate a moment inequality for the sum of $\{-1, 1\}$ random variables (x_1, \dots, x_n) that are k -wise ϵ -biased. We denote their sum by $S = \sum_{i=1}^n x_i$. This will be done by bounding the k^{th} moment of S via the k^{th} moment of the binomial distribution on n uniform independent Bernoulli variables, denoted by B_k . (We assume w.l.o.g. here that k is even). The k^{th} moment of S contains n^k terms, where each term contains at most k variables. More specifically,

$$E[S^k] = E \left[\sum_{i_1 i_2 \dots i_k} x_{i_1} x_{i_2} \dots x_{i_k} \right] = \sum_{i_1 i_2 \dots i_k} E[x_{i_1} x_{i_2} \dots x_{i_k}]$$

Each term in the above summation is of the form $T_i = x_{i_1}^{p_1} x_{i_2}^{p_2} \dots x_{i_r}^{p_r}$, such that $\sum_{j=1}^r p_j = k$. If a term T_i contains a variable whose power is odd, then in B_k its expectation is 0. However, in our case, it follows from the definition of k -wise ϵ -biased random variables (Definitions 2.2 and 2.3), that the expected value of T_i can be at most ϵ . If all the powers in a term T_i are even, then the expected value in both cases is the same and equal to 1. Hence,

Theorem 5.1 *Let S and B_k be as above, then*

$$\text{Prob}[|S| \geq \lambda] \leq \frac{B_k + \epsilon \cdot n^k}{\lambda^k}$$

6 Derandomization

The probability distribution D constructed in Section 3 can be used for derandomizing algorithms. A randomized algorithm \mathcal{A} has a probability space (Ω, P) associated with it, where Ω is the sample space and P is some probability measure. We call a point $w \in \Omega$ a *good* point for some input instance I , if $\mathcal{A}(I, w)$ computes the correct solution. A derandomization of an algorithm means searching the associated sample space Ω for a good point w with respect to a given input instance I . Given such a point w , the algorithm $\mathcal{A}(I, w)$ is now a deterministic algorithm and it is guaranteed to find the correct solution. The problem faced in searching the sample space is that it is generally exponential in size.

In [36, 38, 5] the following strategy was suggested: show that the n probabilistic choices of certain randomized algorithm are only required to be k -wise independent. Hence, a sample space of size $O(n^k)$ suffices. This sample space can be exhaustively searched for a good point

(even in parallel) when k is a constant. A similar strategy can be used with ϵ -biased random variables. First, show that a randomized algorithm has a non-zero probability of success when the probabilistic choices are ϵ -biased. Then, conduct a search of the sample space associated with the variables to find a good point.

This scheme can in principle be applied to all the randomized algorithms for which the limited independence approach was applied. However, we do not necessarily get better algorithms. An attractive feature of our scheme is that random variables that are $\log n$ -wise δ -dependent, for δ which is polynomially small, can be constructed with a *polynomial* sample space. Intuitively, this means that we can achieve “almost” $\log n$ -wise independence with a polynomial sample space (as opposed to $n^{O(\log n)}$).

The k -wise ϵ -biased random variables are especially useful when the proof that a randomized algorithm is successful involves any moment inequality. In that case, one should compute what is the appropriate ϵ such that the error incurred leaves the probability of success non-zero. We exemplify this by showing how a RNC algorithm for the *set balancing* problem can be converted into an NC algorithm.

The second problem we address is finding a heavy codeword in a linear code (Section 6.1). Such a codeword can be found by a simple randomized algorithm. We show how to derandomize it in parallel. The importance of this problem is that it demonstrates that for certain problems it is important to bound the bias of all subsets and not just those of small cardinality.

6.1 Set Balancing

The set balancing problem is defined as follows. A collection of subsets $\mathcal{S} = \{S_1, S_2, \dots, S_n\}$ defined over a base set $B = \{b_1, b_2, \dots, b_n\}$ is given such that the cardinality of each subset is δ . The output is a $\{-1, 1\}$ coloring of B into two sets. Let the 2-coloring of B be denoted by $\mathbf{x} = (x_1, x_2, \dots, x_n)$. The discrepancy of a subset S_i , $\Delta(S_i, \mathbf{x})$, with respect to \mathbf{x} , is defined as $\sum_{j \in S_i} x_j$. The discrepancy of the set family \mathcal{S} is defined as,

$$\Delta(\mathcal{S}, \mathbf{x}) = \max_{S_i} \Delta(S_i, \mathbf{x})$$

Spencer [50, 9] has shown that for each family \mathcal{S} , there exists a 2-coloring \mathbf{x} such that $\Delta(\mathcal{S}, \mathbf{x}) \leq 6\sqrt{n}$. This result is the best possible up to constant factors but it has the drawback of being non-constructive, i.e. does not even imply a probabilistic algorithm. Using the method of conditional probabilities (sequentially), Spencer devised a polynomial-time deterministic algorithm which guarantees a 2-coloring \mathbf{x} such that $\Delta(\mathcal{S}, \mathbf{x}) = O(\sqrt{n \log n})$. This was improved to $O(\sqrt{\delta \log n})$ by Raghavan [45]. For parallel algorithms we cannot guarantee as small a discrepancy as in the sequential case. However, we can come arbitrarily close to the sequential

bounds by computing a τ -good coloring.

Definition 6.1 A 2-coloring \mathbf{x} is τ -good for a set family \mathcal{S} if, for $0 < \tau < \frac{1}{2}$,

$$\Delta(\mathcal{S}, \mathbf{x}) \leq \delta^{0.5+\tau} \sqrt{\log n}$$

There is a simple RNC algorithm to find a τ -good coloring for any \mathcal{S} : pick a random \mathbf{x} uniformly at random from $\{0, 1\}^n$. It turns out that for any \mathcal{S} , the 2-coloring is τ -good with sufficiently high probability. (See for example [34]). This algorithm was derandomized by [13, 39] by proving that $\log n$ -wise independence suffices, and then showing how to conduct a binary search in a sample space of size $n^{O(\log n)}$. (The method of conditional probabilities). We present a direct derandomization which can be implemented in NC^1 .

Assume that B is colored at random; [13, 39] prove the following lemma.

Lemma 6.1 Let $k = a \log n / \log \delta$ be even, where $a > \frac{1}{\tau}$, and let \mathbf{x} be k -wise independent uniform $\{-1, 1\}$ random variables. Then, for any input set family \mathcal{S} , the probability that \mathbf{x} is τ -good is non-zero.

We sketch very briefly the proof idea. Let T_i denote the sum of the random variables in the subset S_i . The main tool used is the k^{th} moment inequality. (See Section 5). More specifically, [13, 39] show that for a given subset S_i ,

$$\text{Prob} \left[\left| T_i - \frac{\delta}{2} \right| > \delta^{0.5+\tau} \sqrt{\log n} \right] \leq \frac{E[T_i^k]}{(\delta^{0.5+\tau} \sqrt{\log n})^k} < \frac{1}{\delta^{\tau k}} < \frac{1}{n}$$

Summing over all subsets, we get that for any set family \mathcal{S} , the probability that a random k -wise coloring is τ -good is non-zero.

What happens when the random variables in \mathbf{x} are k -wise ϵ -biased? We want to choose ϵ such that the probability that the discrepancy in a subset is large, is smaller than $\frac{1}{n}$. Substituting in Theorem 5.1, ϵ and k must be chosen such that

$$\frac{1}{\delta^{\tau k}} + \frac{\epsilon \delta^k}{(\delta^{0.5+\tau} \sqrt{\log n})^k} < \frac{1}{n}$$

We choose $k = \frac{\log 2n}{\tau \log \delta}$, and the above inequality holds if

$$\epsilon < \frac{1}{2n^{1+\frac{1}{\tau}}}$$

Hence, the sample space will be at most of cardinality $n^{O(\frac{1}{\tau})}$ and if τ is a constant, an exhaustive search of the sample space can be conducted. Both the construction and the search can be done in NC^1 .

For the relationship between the set balancing problem and the lattice approximation problem [45], the reader is referred to [39].

6.2 Finding heavy codewords

Let $C[n, k, d]$ be a linear code. In this section we show how to find a codeword whose weight is at least as heavy as the expected weight of a word in C . We call such a codeword, a *heavy* codeword. This problem is a natural generalization of the following problem: given a graph G , find a cut that contains at least half the edges. It is well known that the set of cuts in a graph forms a linear subspace. Let G be the generator matrix of the linear code C . It is easy to see that the expected weight of a codeword in C is $\frac{n}{2}$: if x is chosen uniformly at random, then $E[\text{weight of } Gx] = \frac{n}{2}$.

The latter observation implies a straightforward randomized algorithm for finding a heavy codeword. This algorithm can be made deterministic via the method of conditional probabilities [50, 9] where the entries of x are determined one at a time. How can a heavy codeword be found in parallel? One possible approach for reducing the sample space is by choosing the entries of x to be only k -wise independent (for some $k < n$) and not mutually independent. The difficulty is that the probability of getting a heavy codeword may vanish.

If x is chosen from an ϵ -biased probability space, then $E[\text{weight of } Gx] \geq n(\frac{1-\epsilon}{2})$. If $\epsilon < \frac{1}{2n}$, then there must be a codeword x in the ϵ -biased probability space whose weight is at least $\frac{n}{2}$. This places the problem in NC for the first time. More important, it exhibits that ϵ -biased random variables are also needed as opposed to just k -wise ϵ -biased random variables.

7 Reducing the number of random bits

In this section we present two algorithms for which the number of random bits required can be reduced from linear to logarithmic.

7.1 Matrix multiplication verification

Suppose that three $n \times n$ matrices A, B and C over an arbitrary ring are given; what is the complexity of verifying whether $A \cdot B = C$? Can this be done by avoiding matrix multiplication? This problem was first considered by Freivalds [26] who suggested a randomized algorithm of complexity $O(n^2)$, but it required n random bits. In this section we show how to implement it using only $O(\log n)$ random bits with no time penalty. Our results can readily be generalized to non-square matrices as well.

Let us first review Freivalds' algorithm [26]. Choose a random vector \vec{r} such that each entry in r is picked uniformly at random to be zero or some fixed non-zero element of the ring. Test whether $r \cdot A \cdot B - r \cdot C = 0$. The complexity of applying this procedure is that of multiplying a

matrix by a vector which is obviously $O(n^2)$. For the sake of completeness, we present a proof of his algorithm.

Theorem 7.1 *Suppose that r is a vector chosen in the manner described above. Then with probability at least $\frac{1}{2}$, if $A \cdot B \neq C$, then also $r \cdot A \cdot B \neq r \cdot C$.*

Proof: Let v be a non-zero column vector of $A \cdot B - C$ and let j denote the number of its non-zero entries. The number of different vectors r such that $v \cdot r = 0$ is at most $2^{j-1} \cdot 2^{n-j} = 2^{n-1}$, whereas the number of distinct choices for r is 2^n . Hence, the probability that $v \cdot r = 0$ is at most $\frac{1}{2}$. \square

It follows from the above proof that we can restrict ourselves to the following problem: given a vector v , check whether it is identically zero, where the only operation permitted on the vector is taking its inner product with another vector. Recall from Section 3.1 that a vector r that verifies that a particular vector v is non-zero is called a *distinguisher* (with respect to v). Hence, reducing the number of random bits for the above problem is equivalent to generating a small collection of vectors \mathcal{F} , such that the inner product of a constant fraction of them with *any* non-zero vector is not zero.

The family \mathcal{F} that is constructed in Section 3.1 has this property: with probability at least β , a vector r sampled uniformly will be a distinguisher. The number of random bits needed for sampling \mathcal{F} is $O(\log n)$ and the complexity of the algorithm remains $O(n^2)$.

Notice that if the matrices are over an arbitrary ring, then the construction presented in Section 3.1.2 cannot be used. We can only use the one presented in Section 3.1.1.

7.2 Boolean matrix multiplication

In this section we show how Adi Shamir's boolean matrix multiplication algorithm (see [19, pp. 772-773]) can be implemented with few random bits without increasing the probability of error.

Let $A = (a_{ij})$ and $B = (b_{ij})$ be $n \times n$ Boolean matrices and suppose we would like to multiply them in the quasiring $Q = (\{0, 1\}, \vee, \wedge, 0, 1)$ by using the algorithms for fast matrix multiplications, e.g., Strassen's method. The difficulty is that these methods require that the matrix multiplication be carried out in a ring. This can be handled by working over a large field, but makes multiplication more expensive. Instead, Shamir suggested a randomized algorithm that takes advantage of the fact that the fast methods for matrix multiplication can be carried out in the ring $R = (\{0, 1\}, \oplus, \wedge, 0, 1)$.

We briefly summarize Shamir's algorithm. Let $C = (c_{ij}) = AB$ in the quasiring Q . Generate $A' = (a'_{ij})$ from A using the following randomized procedure:

- If $a_{ij} = 0$, then let $a'_{ij} = 0$.
- If $a_{ij} = 1$, then let $a'_{ij} = 1$ with probability $1/2$ and let $a'_{ij} = 0$ with probability $1/2$. The random choices for each entry are independent.

Let $C' = A'B$ in the ring R . The following lemma is immediate.

Lemma 7.1 *If $c_{ij} = 0$, then $c'_{ij} = 0$. If $c_{ij} = 1$, then $\text{Prob}[c'_{ij} = 1] \geq 1/2$.*

As in the algorithm for verifying matrix multiplication, make the random choices to be ϵ -bias and get that the probability of error is at most $1/2 + \epsilon$.

To decrease the probability of failure, we will run the algorithm for $\log(n^2/\delta)$ choices of the matrix A' . Since the matrix C' has n^2 entries, the total probability of error is bounded from above by δ .

7.3 Simultaneous verification of exponentiation

Suppose that for some prime p and integer a we are given n pairs $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ and we want to verify whether the equality $a^{x_i} = y_i \pmod{p}$ is true for all $1 \leq i \leq n$. A. Fiat and M. Naor (personal communication) have suggested the following randomized algorithm:

1. Pick a random vector $r = (r_1, \dots, r_n)$ where each r_i is chosen uniformly and randomly from $\{0, 1\}$.
2. Compute $t = \sum_{i=1}^n r_i \cdot x_i \pmod{p-1}$ and $m = \prod_{i \in S} y_i^{r_i} \pmod{p}$.
3. Test whether $a^t = m \pmod{p}$.

As in Freivalds' algorithm, if for any $1 \leq i \leq n$, $a^{x_i} \neq y_i$, then the above algorithm will detect it with probability at least $\frac{1}{2}$. The complexity of this algorithm is n multiplications, instead of $n \log p$ for checking each equality separately.

We can actually phrase the problem as that of finding a distinguisher with respect to a non-zero vector in the integer ring modulo $p-1$. Let z_i be such that $a^{z_i} = y_i$. Then, a distinguisher with respect to the vector w must be found, where $w_i = x_i - z_i \pmod{p-1}$. For that we can use the construction of section 3.1 in a similar way to Section 7.1.

Note that the expected size of each entry in Step 4 of the algorithm in Section 3.1.1 is $O(1)$, and therefore, the expected number of multiplications remains $O(n)$.

Suppose that a and N are integers such that $(N, a) = 1$. The above procedure can be applied to verify n equalities of the type $x_i^a = y_i \pmod{N}$. This may be used for instance to check n given RSA equations [46] and thus amortize the cost of verifying signatures.

8 Vector sets for exhaustive testing

A problem that has received much attention in the fault diagnostic literature is that of generating a small set of vectors, $T \subset \{0, 1\}^n$, such that for any k -subset of indices $S = \{i_1, i_2, \dots, i_k\}$, the projection of T on the indices S contains all possible 2^k configurations. See [51] for a bibliography on the problem. Such a set of vectors is called (n, k) -universal. The motivation for this problem is that such a test set allows exhaustive testing of a circuit where each component relies on at most k inputs. Alternatively, one can phrase the problem as searching for a collection of n subsets of a ground set which is as small as possible such that the intersection of any k subsets or their complement is non-empty. This was called k -independent by Kleitman and Spencer [35].

The best known bounds for constructing (n, k) universal sets are given in [51] and [3]. The connection between k -wise δ -dependent probability spaces and small (n, k) -universal sets is made in the next proposition.

Proposition 8.1 *If Ω is a k -wise δ -dependent probability space for $\delta \leq 2^{-k}$, then Ω is also a (n, k) -universal set.*

Proof: If for a k -subset i_1, i_2, \dots, i_k , there is a $\{0, 1\}^k$ configuration which has probability 0 in Ω , then the distance from the uniform distribution of $x_{i_1}, x_{i_2}, \dots, x_{i_k}$ is at least $2^{-(k-1)} > \delta$. \square

Combining the construction suggested in Section 3.1.2 with Lemma 4.2, we can construct a probability space Ω which is k -wise 2^{-k} -dependent such that the cardinality of Ω is $\log n \cdot 2^{O(k)}$. Thus, our construction for a (n, k) universal set can be phrased in coding theory terminology. Let G be the generating matrix of a binary $[n', k', d]$ linear code for which

$$\frac{d}{n'} = \frac{1}{2} - \frac{1}{2^{2k}}$$

and let H be the parity check matrix of a binary $[n, n', k]$ linear code. The (n, k) universal set consists of the rows of the matrix $G \cdot H$. This is better than the best explicit construction given in [51] for $k \ll n$, and matches the lower bound given there up to constant factors. Alon [3] has shown an explicit construction of size $\log n \cdot 2^{O(k^2)}$ which is optimal for constant k . No such constructions were known for larger values of k . For k which is $\Theta(\log n)$, this is the first explicit construction of an (n, k) -universal set of polynomial size.

9 Communication complexity

Suppose player A has a string $x \in \{0, 1\}^n$ and player B has a string $y \in \{0, 1\}^n$, and they wish to decide whether $x = y$ while exchanging as few bits as possible. It is well known that this can be done by exchanging $O(\log n)$ bits and the probability of the protocol arriving at the correct

result is at least a constant. We can show how to achieve probability of success which is any $\frac{1}{\text{poly}(n)}$, while maintaining the logarithmic communication complexity.

The algorithm will use the first two stages in constructing the small biased probability spaces. To achieve probability of error ϵ :

1. Player A chooses $O(\log n + \log \frac{1}{\epsilon})$ random bits that define vectors r_1, r_2, \dots, r_l . She sends the random bits to player B and also sends $\langle r_1, x \rangle, \langle r_2, x \rangle, \dots, \langle r_l, x \rangle$.
2. B computes $\langle r_1, y \rangle, \langle r_2, y \rangle, \dots, \langle r_l, y \rangle$. If for any i , $\langle r_i, y \rangle \neq \langle r_i, x \rangle$ he announces it. Otherwise, he concludes that $x = y$.

If $x \neq y$, then with probability at least $1 - \epsilon$, for at least one i , $\langle r_i, x \oplus y \rangle \neq 0$ and hence $\langle r_i, x \rangle \neq \langle r_i, y \rangle$.

Once a distinguisher has been found, detecting an index of an entry on which players A and B differ is easy by exchanging $\log n$ bits.

Another application is for the problem of determining set equality. Suppose that $A, B \subset \{1, \dots, n\}$ and we wish to decide whether $A = B$. We would like a single pass on the elements of A and B and the amount of memory should be $O(\log n)$. Here again a method that achieves constant probability error is known (see Blum and Kannan [16] for a description).

We will show that it is possible to achieve any ϵ error while using only $O(\log n + \log \frac{1}{\epsilon})$ bits of memory. Let v_A and v_B denote the incidence vectors of A and B . We can think of the problem as deciding whether $v_A \oplus v_B = 0$. Again, we can use the first two stages of the construction of small bias probability spaces.

1. Pick $O(\log n + \log \frac{1}{\epsilon})$ random bits that define r_1, r_2, \dots, r_l . (They are not computed explicitly at this point.) Let $r_i(a)$ denote the a^{th} coordinate of r_i .
2. Initialize bits d_1, d_2, \dots, d_l to 0.
3. For each element $a \in A$ and for each $1 \leq i \leq l$, $d_i \leftarrow d_i \oplus r_i(a)$.
4. For each element $b \in B$ and for each $1 \leq i \leq l$, $d_i \leftarrow d_i \oplus r_i(b)$.
5. If all the d_i -s are 0, decide that $A = B$, otherwise $A \neq B$.

As before, if for at least one i , $\langle r_i, v_A \oplus v_B \rangle \neq 0$, we will detect the inequality of the sets. The probability that this happens is at least $1 - \epsilon$. For this application, the first method of constructing small bias probability spaces should be used, since we are not interested in the full vector r_i , but in selected locations of it.

This can be looked upon as a family \mathcal{H} of hash functions with the following property. Let $h \in \mathcal{H}$, where $h : \{1 \dots n\} \rightarrow \{1 \dots m\}$. The family \mathcal{H} is accessible with $O(\log n + \log m)$ bits and for any subsets A and B of $\{1 \dots n\}$, the probability that $\sum_{a \in A} h(a) = \sum_{b \in B} h(b)$ is smaller than $\frac{1}{m^\beta}$ for β constant, where addition is bitwise XOR.

10 Further results

Since the preliminary version of this paper appeared in [42], several new applications of small bias probability spaces have been discovered. Alon [4] used them to obtain an NC¹ algorithm for the parallel version of Beck's algorithm for the Lovasz local lemma. Feder, Kushilevitz and Naor [24] have applied them to amortize the communication complexity of equality. Blum et al. [15] used the hash function construction to authenticate memories. Kushilevitz and Mansour [31] used small bias probability spaces to derandomize their decision tree learning algorithm. Alon et al. [7] have used the polylogarithmic size construction of $\log \log n$ -wise $1/\log n$ -bias probability space in order to derandomize an algorithm for all pairs shortest path whose running time is almost that of matrix multiplication. Blum and Rudich [14] have used the construction of (n, k) -universal sets of Section 8 in order to derandomize a k -term DNF polynomial time learning algorithm, for k which is logarithmic in the number of variables.

Azar, Motwani and Naor [10] defined and constructed small bias probability spaces for non-binary random variables. Even et al. [23] constructed δ -dependent non-uniform probability spaces based on small bias probability spaces. Alon et al. [8] provided simple and different constructions for small bias probability spaces. Their constructions are based on quadratic characters and linear feedback shift registers. Schulman [48] constructed efficiently probability spaces with known dependencies.

Acknowledgement

We would like to thank Shuki Bruck for his contribution to Section 3.1.2. We would also like to thank Noga Alon, Yishay Mansour and Avi Wigderson for helpful discussions. We thank Steven Ponzio and the referees for their careful reading of the manuscript. We thank Robert Beals for pointing out [22].

References

- [1] M. Ajtai, J. Komlos and E. Szemerédi, Deterministic simulation in LOGSPACE, Proceedings of the 19th ACM Annual Symposium on Theory of Computing, (1987) pp. 132-140.

- [2] M. Ajtai and A. Wigderson, Deterministic simulation of probabilistic constant depth circuits, Proceedings of the 26th IEEE Symposium on Foundations of Computer Science (1985), pp. 11-19.
- [3] N. Alon, Explicit constructions of exponential sized families of k -independent sets, Discrete Math, 58, pp. 191-193 (1986).
- [4] N. Alon, A parallel algorithmic version of the local lemma, Random Structures and Algorithms, Vol. 2 (1991), pp. 367-378.
- [5] N. Alon, L. Babai and A. Itai, A fast and simple randomized parallel algorithm for the maximal independent set problem, Journal of Algorithms, Vol 7 (1986) , pp. 567-583.
- [6] N. Alon, J. Bruck, J. Naor, M. Naor and R. Roth, Construction of Asymptotically Good Low-Rate Error-Correcting Codes through Pseudo-Random Graphs, IEEE Transactions on Information Theory, vol. 38 (1992), pp. 509-516.
- [7] N. Alon, Z. Galil, O. Margalit and M. Naor, Witnesses for matrix multiplication and for shortest paths, Proceedings of the 33rd IEEE Symposium on Foundations of Computer Science, (1992) pp. 417-426.
- [8] N. Alon, O. Goldreich, J. Hastad and R. Peralta, Simple constructions for almost k -wise independent random variables, Random Structures and Algorithms, Vol. 3 (1992), pp. 289-304.
- [9] N. Alon and J. Spencer, **The probabilistic method**, John Wiley and Sons Inc., New York, 1992.
- [10] Y. Azar, R. Motwani and J. Naor Approximating arbitrary distributions using small sample spaces, manuscript, 1990.
- [11] E. Bach, Realistic analysis of some randomized algorithms, Proceedings of the 19th annual ACM Symposium on Theory of Computing, (1987), pp. 453-461.
- [12] R. Ben-Nathan, M.Sc. Thesis, Hebrew University (1990).
- [13] B. Berger and J. Rompel, Simulating $(\log^c n)$ -wise independence in NC, Journal of the ACM, Vol. 38 (1991), pp. 1026-1046.
- [14] A. Blum and S. Rudich, Fast learning k -term DNF formulas with queries, To appear, Proceedings of the 24th annual ACM Symposium on Theory of Computing, (1992).
- [15] M. Blum, W. Evans, P. Gemmell, S. Kannan and M. Naor, Checking the correctness of memories, Proceedings of the 32nd IEEE Symposium on Foundations of Computer Science, (1991), pp. 90-99.

- [16] M. Blum and R. Kannan, Designing programs that check their work, Proceedings of the 21st annual ACM Symposium on Theory of Computing, (1987) pp. 86-97.
- [17] B. Chor and O. Goldreich, On the power of two-point based sampling, Journal of Complexity, Vol 5 (1989), pp. 96-106.
- [18] B. Chor, O. Goldreich, J. Hastad, J. Friedman, S. Rudich and R. Smolensky, The bit extraction problem or t -resilient functions, Proceedings of the 26th IEEE Symposium on Foundations of Computer Science, (1985) pp. 396-407.
- [19] T. H. Cormen, C. E. Leiserson and R. L. Rivest, Introduction to Algorithms, The MIT Press, Cambridge, Massachusetts, 1991.
- [20] A. Cohen and A. Wigderson, Dispersers, deterministic amplification and weak random sources, Proceedings of the 30th IEEE Symposium on Foundations of Computer Science, (1989) pp. 14-19.
- [21] A. Cohen and A. Wigderson, Multigraph Amplification, Survey (1989).
- [22] P. Diaconis, Group Representations in Probability and Statistics, IMS Lecture Notes-Monograph Series, Volume 11, Hayward, CA (1988).
- [23] G. Even, O. Goldreich, M. Luby, N. Nisan and B. Velickovic, Approximation of general independent distributions, , Proceedings of the 24th annual ACM Symposium on Theory of Computing, (1992), pp. 10-16.
- [24] T. Feder, E. Kushilevitz and M. Naor, Amortized communication complexity, Proceedings of the 32nd IEEE Symposium on Foundations of Computer Science, (1991), pp. 239-248.
- [25] W. Feller, **An Introduction to probability theory and its applications**, John Wiley, 1968.
- [26] R. Freivalds, Fast probabilistic algorithms, Springer Verlag Lecture Notes in CS #74, Mathematical Foundations of CS, pp. 57-69 (1979).
- [27] O. Gaber and Z. Galil, Explicit construction of linear size superconcentrators, Journal Computer Systems Sciences, 22, (1981) pp. 407-420.
- [28] R. Impagliazzo and D. Zuckerman, Recycling random bits, Proceedings of the 30th IEEE Symposium on Foundations of Computer Science, (1989) pp. 248-253.
- [29] S. Jimbo and A. Marouka, Expanders obtained from affine transformations, Proceedings of the 17th annual ACM Symposium on Theory of Computing, (1985) pp. 88-97.

- [30] J. Justesen, A class of asymptotically good algebraic codes, *IEEE Transactions on Information Theory*, Vol. 18 (1972), pp. 652-656.
- [31] E. Kushilevitz and Y. Mansour, Learning decision trees using the Fourier spectrum, *Proceedings of the 23rd annual ACM Symposium on Theory of Computing*, (1991) pp. 455-464.
- [32] R. Karp and N. Pippenger, A time randomness tradeoff, *AMS conference on probabilistic computation and complexity*, Durham NC, (1983).
- [33] H. Karloff and P. Raghavan, Randomized algorithms and pseudorandom numbers, *Proceedings of the 20th annual ACM Symposium on Theory of Computing*, (1988), pp. 310-321.
- [34] H. Karloff and D. Shmoys, Efficient parallel algorithms for edge coloring problems, *J. of Algorithms*, Vol. 8 (1987), pp. 39-52.
- [35] D. J. Kleitman and J. Spencer, Families of k -independent sets, *Discrete Math*, Vol 6 (1973), pp. 255-262.
- [36] R. M. Karp and A. Wigderson, A fast parallel algorithm for the maximal independent set problem, *Journal of the ACM*, vol. 32 (1985), pp. 762-773.
- [37] A. Lubotzky, R. Phillips and P. Sarnak, Explicit expanders and the Ramanujan conjecture, *Proceedings of the 18th annual ACM Symposium on Theory of Computing*, (1986), pp. 240-246.
- [38] M. Luby, A simple parallel algorithm for the maximal independent set problem, *Siam J. Comp.*, Vol 15 (1986), pp. 1036-1053.
- [39] R. Motwani, J. Naor and M. Naor, The probabilistic method yields deterministic parallel algorithms, *Proceedings of the 30th IEEE Symposium on Foundations of Computer Science*, (1989), pp. 8-13.
- [40] F. J. MacWilliams and N. J. A. Sloane, *The theory of error correcting codes*, North Holland, Amsterdam, 1977.
- [41] M. Naor, *Constructing Ramsey Graphs from Small Probability Spaces*, IBM Research Report RJ 8810, 1992.
- [42] J. Naor and M. Naor, Small-bias probability spaces: efficient constructions and applications, *Proceedings of the 22nd annual ACM Symposium on Theory of Computing* (1990), pp. 213-223.
- [43] N. Nisan and A. Wigderson, Hardness vs. Randomness, *Proceedings of the 29th IEEE Symposium on Foundations of Computer Science* (1988), pp. 2-11.

- [44] R. Peralta, On the randomness complexity of algorithms, University of Wisconsin, Milwaukee, CS Research Report TR 90-1.
- [45] P. Raghavan, Probabilistic construction of deterministic algorithms: approximating packing integer programs, *Journal of Computer Systems and Science*, Vol. 37 (1988), pp. 130-143.
- [46] R. Rivest, A. Shamir and L. Adelman, A method for obtaining digital signatures and public key cryptosystems, *Communications of the ACM*, vol 21 (1978) pp. 120-126.
- [47] M. Santha, On using deterministic functions to reduce randomness in probabilistic algorithms, *Information and Computation*, Vol. 74 (1987), pp. 241-249.
- [48] L. Schulman, Sample spaces uniform on neighborhoods, *Proceedings of the 24th annual ACM Symposium on Theory of Computing*, (1992), pp. 17-25.
- [49] M. Sipser, Expanders, randomness, or time versus space, *Journal of Computer Systems Sciences*, Vol. 36, (1988) pp. 379-383.
- [50] J. Spencer, **Ten lectures on the probabilistic method**, SIAM (Philadelphia), 1987.
- [51] G. Seroussi and N. Bshouti, Vector sets for exhaustive testing of logic circuits, *IEEE Transactions on Information Theory*, Vol. 34, pp. 513-522 (1988).
- [52] U. Vazirani, Randomness, adversaries and computation, Ph.D. Thesis, University of California, Berkeley (1986).