

Balls and Bins: Smaller Hash Families and Faster Evaluation

L. Elisa Celis
ecelis@cs.washington.edu
University of Washington

Omer Reingold
omer.reingold@microsoft.com
Microsoft Research SVC

Gil Segev
gil.segev@microsoft.com
Microsoft Research SVC

Udi Wieder
uwieder@microsoft.com
Microsoft Research SVC

Abstract— A fundamental fact in the analysis of randomized algorithms is that when n balls are hashed into n bins independently and uniformly at random, with high probability each bin contains at most $O(\log n / \log \log n)$ balls. In various applications, however, the assumption that a truly random hash function is available is not always valid, and explicit functions are required.

In this paper we study the size of families (or, equivalently, the description length of their functions) that guarantee a maximal load of $O(\log n / \log \log n)$ with high probability, as well as the evaluation time of their functions. Whereas such functions must be described using $\Omega(\log n)$ bits, the best upper bound was formerly $O(\log^2 n / \log \log n)$ bits, which is attained by $O(\log n / \log \log n)$ -wise independent functions. Traditional constructions of the latter offer an evaluation time of $O(\log n / \log \log n)$, which according to Siegel’s lower bound [FOCS ’89] can be reduced only at the cost of significantly increasing the description length.

We construct two families that guarantee a maximal load of $O(\log n / \log \log n)$ with high probability. Our constructions are based on two different approaches, and exhibit different trade-offs between the description length and the evaluation time. The first construction shows that $O(\log n / \log \log n)$ -wise independence can in fact be replaced by “gradually increasing independence”, resulting in functions that are described using $O(\log n \log \log n)$ bits and evaluated in time $O(\log n \log \log n)$. The second construction is based on derandomization techniques for space-bounded computations combined with a tailored construction of a pseudorandom generator, resulting in functions that are described using $O(\log^{3/2} n)$ bits and evaluated in time $O(\sqrt{\log n})$. The latter can be compared to Siegel’s lower bound stating that $O(\log n / \log \log n)$ -wise independent functions that are evaluated in time $O(\sqrt{\log n})$ must be described using $\Omega(2^{\sqrt{\log n}})$ bits.

1. INTRODUCTION

Traditional analysis of randomized algorithms and data structures often assumes the availability of a truly random function, whose description length and evaluation time are not taken into account as part of the overall performance. In various applications, however, such an assumption is not valid and explicit constructions are required. This motivated a well-studied line of research aiming at designing explicit and rather small families of functions, dating back more than 30 years to the seminal work of Carter and Wegman [5].

In this paper we study explicit constructions of families for the classical setting of hashing n balls into n bins. A well-known and useful fact is that when n balls are hashed into n bins independently and uniformly at random, with high probability each bin contains at most

$O(\log n / \log \log n)$ balls. Thus, a natural problem is to construct explicit and significantly smaller families of functions that offer the same maximal load guarantee. More specifically, we are interested in families \mathcal{H} of functions that map a universe U into the set $\{1, \dots, n\}$, such that for any set $S \subseteq U$ of size n a randomly chosen function $h \in \mathcal{H}$ guarantees a maximal load of $O(\log n / \log \log n)$ with high probability. The main measures of efficiency for such families are the description length and evaluation time of their functions.

It is well-known that any family of $O(\log n / \log \log n)$ -wise independent functions guarantees a maximal load of $O(\log n / \log \log n)$ with high probability, and this already yields a significant improvement over a truly random function. Specifically, such functions can be represented by $O(\log^2 n / \log \log n)$ bits, instead of $O(|U| \log n)$ bits for a truly random function¹. A natural approach for reducing the description length is to rely on k -wise independence for $k = o(\log n / \log \log n)$, but so far no progress has been made in this direction (even though to the best of our knowledge an explicit lower bound is only known for $k = 2$ [1]). At the same time, a standard application of the probabilistic method shows that there exists such a family where each function is described by only $O(\log n)$ bits, which is in fact optimal. This naturally leads to the following open problem (whose variants were posed explicitly by Alon et al. [1] and by Pagh et al. [26]):

Problem 1: Construct an explicit family that guarantees a maximal load of $O(\log n / \log \log n)$ with high probability, in which each function is described by $o(\log^2 n / \log \log n)$ bits, or even $O(\log n)$ bits.

In terms of the evaluation time, an $O(\log n / \log \log n)$ -wise independent function can be evaluated using traditional constructions in time $O(\log n / \log \log n)$. A lower bound proved by Siegel [30] shows that the latter can be reduced only at the cost of significantly increasing the description length. For example, even for $k = O(\log n / \log \log n)$ a constant evaluation time requires polynomial space. In the same work Siegel showed a tight (but impractical) matching upper bound. Subsequent constructions improve the constants considerably (see Section 1.2 for a more elaborated

¹For simplicity we assume $|U|$ is polynomial in n , as otherwise one can reduce the size of U using a pair-wise independent function.

discussion), but suffer from descriptions of length at least n^ϵ bits for a small constant $\epsilon > 0$. This leads to the following open problem:

Problem 2: Construct an explicit family that guarantees a maximal load of $O(\log n / \log \log n)$ with high probability, in which each function is evaluated in time $o(\log n / \log \log n)$ and represented by $n^{o(1)}$ bits.

1.1. Our Contributions

We present two constructions of hash families that guarantee a maximal load of $O(\log n / \log \log n)$ when hashing n elements into n bins with all but an arbitrary polynomially-small probability. These are the first explicit constructions in which each function is described using less than $O(\log^2 n / \log \log n)$ bits. Our constructions offer different trade-offs between the description length of the functions and their evaluation time. Table I summarizes the parameters of our constructions with prior work.

Construction 1 – gradually-increasing independence:

In our first construction each function is described using $O(\log n \log \log n)$ bits and evaluated in time $O(\log n \log \log n)$. Whereas $O(\log n / \log \log n)$ -wise independence suffices for a maximal load of $O(\log n / \log \log n)$, the main idea underlying our construction is that *the entire output* need not be $O(\log n / \log \log n)$ -wise independent.

Our construction is based on concatenating the outputs of $O(\log \log n)$ functions which are *gradually* more independent: each function f in our construction is described using d functions h_1, \dots, h_d , and for any $x \in [u]$ we define

$$f(x) = h_1(x) \circ \dots \circ h_d(x) ,$$

where we view the output of each h_i as a binary string, and \circ denotes the concatenation operator on binary strings. The first function h_1 is only $O(1)$ -wise independent, and the level of independence gradually increases to $O(\log n / \log \log n)$ -wise independence for the last function h_d . As we increase the level of independence, we decrease the output length of the functions from $\Omega(\log n)$ bits for h_1 to $O(\log \log n)$ bits for h_d . We instantiate these $O(\log \log n)$ functions using ϵ -biased distributions. The trade-off between the level of independence and the output length implies that each of these functions can be described using only $O(\log n)$ bits and evaluated in time $O(\log n)$.

Construction 2 – derandomizing space-bounded computations: In our second construction each function is described using $O(\log^{3/2} n)$ bits and evaluated in time $O(\log^{1/2} n)$. Each function f in our construction is described using a function h that is $O(1)$ -wise independent, and $\ell = O(2^{\log^{1/2} n})$ functions g_1, \dots, g_ℓ that are $O(\log^{1/2} n)$ -wise independent, and for any $x \in [u]$ we define

$$f(x) = g_{h(x)}(x) .$$

Naively, the description length of such a function f is $O(\ell \cdot \log^{3/2} n)$ bits, and the main idea underlying our

approach is that instead of sampling the functions g_1, \dots, g_ℓ independently and uniformly at random, they can be obtained as the output of an explicit pseudorandom generator for space-bounded computations using a seed of length $O(\log^{3/2} n)$ bits. Moreover, we present a new construction of a pseudorandom generator for space-bounded computations in which the description of each of these ℓ functions can be computed in time $O(\log^{1/2} n)$ without increasing the length of the seed.

Our generator is obtained as a composition of those constructed by Nisan [23] and by Nisan and Zuckerman [24] together with an appropriate construction of a randomness extractor for instantiating the Nisan-Zuckerman generator. The evaluation time of our second construction can be compared to Siegel’s lower bound [30] stating that $O(\log n / \log \log n)$ -wise independent functions that are evaluated in time $O(\log^{1/2} n)$ must be described using $\Omega(2^{\log^{1/2} n})$ bits.

We note that a generator with an optimal seed length against space-bounded computations will directly yield a hash family with the optimal description length $O(\log n)$ bits. Unfortunately, the best known generator [23] essentially does not give any improvement over using $O(\log n / \log \log n)$ -wise independence. Instead, our above-mentioned approach is based on techniques that were developed in the area of pseudorandomness for space-bounded computations with which we obtain an improvement in our specific setting. Specifically, our construction is inspired by the generator constructed by Lu [19] for the simpler class of combinatorial rectangles.

Extensions: It is possible to show that the hash families constructed in this paper can be successfully employed for storing elements using linear probing. In this setting our constructions guarantee an insertion time of $O(\log n)$ with high probability when storing $(1 - \alpha)n$ elements in a table of size n , for any constant $0 < \alpha < 1$ (and have constant expected insertion time as follows from [26]). Prior to our work, constructions that offered such a high probability bound had either description length of $\Omega(\log^2 n)$ bits with $\Omega(\log n)$ evaluation time (using $O(\log n)$ -wise independence [29]) or description length of $\Omega(n^\epsilon)$ bits with constant evaluation time [30], [28].

In addition, our constructions can easily be augmented to offer $O(\log \log n)$ -wise independence (for the first construction), and $O(\log^{1/2} n)$ -wise independence (for the second construction) without affecting their description length and evaluation time. This may be useful, for example, in any application that involves tail bounds for limited independence.

Lower bounds: We accompany our constructions with two somewhat folklore lower bounds. First, for a universe of size at least n^2 , any family of functions has a maximal load of $\Omega(\log n / \log \log n)$ with high probability. Second, the functions of any family that guarantees a maximal load of $O(\log n / \log \log n)$ with probability $1 - \epsilon$ must be described

| | Description length (bits) | Evaluation time |
|---|--|--|
| Simulating full independence ([11], [25]) | $O(n \log n)$ | $O(1)$ |
| [30],[8],[28] | n^ϵ (for constant $\epsilon < 1$) | $O(1)$ |
| $O\left(\frac{\log n}{\log \log n}\right)$ -wise independence (polynomials) | $O\left(\frac{\log^2 n}{\log \log n}\right)$ | $O\left(\frac{\log n}{\log \log n}\right)$ |
| This paper (Section 4) | $O\left(\log^{3/2} n\right)$ | $O\left(\log^{1/2} n\right)$ |
| This paper (Section 3) | $O(\log n \log \log n)$ | $O(\log n \log \log n)$ |

Table I
THE DESCRIPTION LENGTH AND EVALUATION TIME OF OUR CONSTRUCTIONS AND PREVIOUSLY KNOWN CONSTRUCTIONS THAT GUARANTEE A MAXIMAL LOAD OF $O(\log n / \log \log n)$ W.H.P.

by $\Omega(\log n + \log(1/\epsilon))$ bits. Due to space constraints, the proofs of these lower bounds are contained in the full version of this paper [6].

1.2. Related Work

As previously mentioned, a truly random function guarantees a maximal load of $O(\log n / \log \log n)$ with high probability, but is described by $\Omega(u \log n)$ bits. Pagh and Pagh [25] and Dietzfelbinger and Woelfel [11], in a rather surprising and useful result, showed it is possible to simulate full independence for any specific set of size n (with high probability) using only $O(n \log n)$ bits and constant evaluation time. A different and arguably simpler construction was later proposed by Dietzfelbinger and Rink [10].

In an influential work, Siegel [30] showed that for a small enough constant $\epsilon > 0$ it is possible to construct a family of functions where there is a small probability of error, but if error is avoided then the family is n^ϵ -wise independent, and each function is described using $n^{\epsilon'}$ bits (where $\epsilon < \epsilon' < 1$). More importantly, a function is evaluated in constant time. While this construction has attractive asymptotic behavior it seems somewhat impractical, and was improved by Dietzfelbinger and Rink [10] who proposed a more practical construction (offering the same parameters). Siegel [30] also proved a cell probe time-space tradeoff for computing almost k -wise independent functions. Assuming that in one time unit we can read a word of $\log n$ bits, denote by Z the number of words in the representation of the function and by T the number of probes to the representation. Siegel showed that when computing a k -wise δ -dependent function into $[n]$ then either $T \geq k$ or $Z \geq n^{\frac{1}{k}}(1 - \delta)$. Observe that if k is a constant, setting $T \leq k - 1$ already implies the space is polynomial in n . Also, computing a $O(\log n)$ -wise independent function in time $O(\sqrt{\log n})$ requires the space to be roughly $O(2^{\sqrt{\log n}})$.

Few constructions diverged from the large k -wise independence approach. In [1] it is shown that matrix multiplication over \mathbb{Z}_2 yields a maximal load of $O(\log n \log \log n)$ with high probability, where each function is described using $O(\log^2 n)$ bits and evaluated in time $O(\log n)$. Note that this

family is only pairwise independent. The family of functions described in [8] (which is $O(1)$ -wise independent) yields a maximal load of $O(\log n / \log \log n)$ with high probability, where each function is described using n^ϵ bits and evaluated in constant time (similar to [30]). The main advantage of this family is its practicality: it is very simple and the constants involved are small.

Recently, Pătraşcu and Thorup [28] showed a another practical and simple construction that uses n^ϵ space and $O(1)$ time and can replace truly random functions in various applications, although it is only 3-wise independent. A different approach was suggested by Mitzenmacher and Vadhan [21], who showed that in many cases a pair-wise independent function suffices, provided the hashed elements themselves have a certain amount of entropy.

1.3. Paper Organization

In Section 2 we present the basic definitions and tools used in our constructions. In Sections 3 and 4 we present our first and second constructions, respectively. We conclude in Section 5 with extensions and open problems. Please see [6] for a full version of this paper.

2. PRELIMINARIES AND TOOLS

We now present the background used in our constructions.

2.1. Basic Definitions and the Computational Model

Throughout this paper, we consider \log to be of base 2. For an integer $n \in \mathbb{N}$ we denote by $[n]$ the set $\{1, \dots, n\}$, and by U_n the uniform distribution over the set $\{0, 1\}^n$. For a random variable X we denote by $x \leftarrow X$ the process of sampling a value x according to the distribution of X . Similarly, for a finite set S we denote by $x \leftarrow S$ the process of sampling a value x according to the uniform distribution over S . The *statistical distance* between two random variables X and Y over a finite domain Ω is $\text{SD}(X, Y) = \frac{1}{2} \sum_{\omega \in \Omega} |\Pr[X = \omega] - \Pr[Y = \omega]|$. For two bit-strings x and y we denote by $x \circ y$ their concatenation.

We consider the *unit cost RAM model* in which the elements are taken from a universe of size u , and each element can be stored in a single word of length $w = O(\log u)$ bits. Any operation in the standard instruction set can be executed in constant time on w -bit operands. This includes addition, subtraction, bitwise Boolean operations, parity, left and right bit shifts by an arbitrary number of positions, and multiplication. The unit cost RAM model has been the subject of much research, and is considered the standard model for analyzing the efficiency of data structures and hashing schemes (see, for example, [9], [15], [16], [20], [25] and the references therein).

2.2. Random Variables and Limited Independence

A family \mathcal{F} of functions $f : [u] \rightarrow [v]$ is *k -wise δ -dependent* if for any distinct $x_1, \dots, x_k \in [u]$ the statistical

distance between the distribution $(f(x_1), \dots, f(x_k))$ where $f \leftarrow \mathcal{F}$ and the uniform distribution over $[v]^k$ is at most δ . A simple example for k -wise independent functions (with $\delta = 0$) is the family of all polynomials of degree $k - 1$ over a finite field. Each polynomial can be represented using $O(k \max\{\log u, \log v\})$ bits and evaluated in time $O(k)$ assuming field elements fit into a constant number of words.

For our constructions we require functions that have a more succinct representation, and still enjoy a fast evaluation. For this purpose we implement k -wise δ -dependent functions using ϵ -biased distributions [22]. A sequence of random variables X_1, \dots, X_n over $\{0, 1\}$ is ϵ -biased if for any non-empty set $S \subseteq [n]$ it holds that $|\Pr[\bigoplus_{i \in S} X_i = 1] - \Pr[\bigoplus_{i \in S} X_i = 0]| \leq \epsilon$, where \oplus is the exclusive-or operator on bits. Alon et al. [2, Sec. 5] constructed an ϵ -biased distribution over $\{0, 1\}^n$ where each point $x \in \{0, 1\}^n$ in the sample space can be specified using $O(\log(n/\epsilon))$ bits, and each individual bit of x can be computed in time $O(\log(n/\epsilon))$. Moreover, in the unit cost RAM model with a word size of $w = \Omega(\log(n/\epsilon))$ bits, each block of $t \in [n]$ consecutive bits can be computed in time $O(\log(n/\epsilon) + t)$.²

Using the fact that for any k , an ϵ -biased distribution is also k -wise δ -dependent for $\delta = \epsilon^{2k/2}$ (see, for example, [2, Cor. 1]), we obtain the following corollary:

Corollary 2.1. *For any integers u and v such that v is a power of 2, there exists a family of k -wise δ -dependent functions $f : [u] \rightarrow [v]$ where each function can be described using $O(\log u + k \log v + \log(1/\delta))$ bits. Moreover, in the unit cost RAM model with a word size of $w = \Omega(\log u + k \log v + \log(1/\delta))$ bits each function can be evaluated in time $O(\log u + k \log v + \log(1/\delta))$.*

The construction is obtained from the ϵ -biased distribution of Alon et al. over $n = u \log v$ bits with $\epsilon = \delta^{2^{-k} \log v / 2}$. One partitions the $u \log v$ bits into u consecutive blocks of $\log v$ bits, each of which represents a single output value in the set $[v]$.

A useful tail bound for limited independence: The following is a natural generalization of a well-known tail bound for $2k$ -wise independent random variables [4, Lemma 2.2] (see also [12]) to random variables that are $2k$ -wise δ -dependent. See the full version of this paper [6] for a proof.

Lemma 2.2. *Let $X_1, \dots, X_n \in \{0, 1\}$ be $2k$ -wise δ -dependent random variables, for some $k \in \mathbb{N}$ and $0 \leq \delta < 1$, and let $X = \sum_{i=1}^n X_i$ and $\mu = \mathbb{E}[X]$. Then, for any $t > 0$ it holds that*

$$\Pr[|X - \mu| > t] \leq 2 \left(\frac{2nk}{t^2} \right)^k + \delta \left(\frac{n}{t} \right)^{2k}.$$

²Specifically, the seed consists of two elements $x, y \in \text{GF}[2^m]$, where $m = O(\log(n/\epsilon))$, and the i th output bit is the inner product (modulo 2) of the binary representations of x^i and y .

2.3. Randomness Extraction

The *min-entropy* of a random variable X is $H_\infty(X) = -\log(\max_x \Pr[X = x])$. A k -source is a random variable X with $H_\infty(X) \geq k$. A (T, k) -block source is a random variable $X = (X_1, \dots, X_T)$ where for every $i \in [T]$ and x_1, \dots, x_{i-1} it holds that $H_\infty(X_i | X_1 = x_1, \dots, X_{i-1} = x_{i-1}) \geq k$. In our setting we find it convenient to rely on the following natural generalization of block sources:

Definition 2.3. *A random variable $X = (X_1, \dots, X_T)$ is a (T, k, ϵ) -block source if for every $i \in [T]$ it holds that if $(x_1, \dots, x_{i-1}) \leftarrow (X_1, \dots, X_{i-1})$, then*

$$\Pr[H_\infty(X_i | X_1 = x_1, \dots, X_{i-1} = x_{i-1}) \geq k] \geq 1 - \epsilon.$$

The following lemma and corollary show that any source with high min-entropy can be viewed as a (T, k, ϵ) -block source.

Lemma 2.4 ([14]). *Let X_1 and X_2 be random variables over $\{0, 1\}^{n_1}$ and $\{0, 1\}^{n_2}$, respectively, such that $H_\infty(X_1 X_2) \geq n_1 + n_2 - \Delta$. Then, $H_\infty(X_1) \geq n_1 - \Delta$, and for any $\epsilon > 0$ it holds that*

$$\Pr_{x_1 \leftarrow X_1} [H_\infty(X_2 | X_1 = x_1) < n_2 - \Delta - \log(1/\epsilon)] < \epsilon.$$

Corollary 2.5. *Any random variable $X = (X_1, \dots, X_T)$ over $(\{0, 1\}^n)^T$ with $H_\infty(X) \geq Tn - \Delta$ is a $(T, n - \Delta - \log(1/\epsilon), \epsilon)$ -block source for any $\epsilon > 0$.*

The following defines the notion of a strong randomness extractor.

Definition 2.6. *A function $\text{Ext} : \{0, 1\}^n \times \{0, 1\}^d \rightarrow \{0, 1\}^m$ is a strong (k, ϵ) -extractor if for any k -source X over $\{0, 1\}^n$ it holds that*

$$\text{SD}((S, \text{Ext}(X, S)), (S, Y)) \leq \epsilon,$$

where S and Y are independently and uniformly distributed over $\{0, 1\}^d$ and $\{0, 1\}^m$, respectively.

For our application we rely on the generalization of the leftover hash lemma to block sources [7], [17], [33], showing that a strong extractor enables to reuse the same seed for a block source. This generalization naturally extends to (T, k, ϵ) -block sources:

Lemma 2.7. *Let $X = (X_1, \dots, X_T)$ be a (T, k, ϵ) -block source over $\{0, 1\}^n$ and let \mathcal{H} be a family of pairwise independent functions $h : \{0, 1\}^n \rightarrow \{0, 1\}^m$, where $m \leq k - 2 \log(1/\epsilon)$. Then,*

$$\text{SD}((h, h(X_1), \dots, h(X_T)), (h, Y_1, \dots, Y_T)) \leq 2T\epsilon,$$

where $h \leftarrow \mathcal{H}$, and (Y_1, \dots, Y_T) is independently and uniformly distributed over $(\{0, 1\}^m)^T$.

2.4. Pseudorandom Generators for Space-Bounded Computations

In this paper we model space-bounded computations as layered branching programs (LBP)³. An (s, v, ℓ) -LBP is a directed graph with $2^s(\ell + 1)$ vertices that are partitioned into $\ell + 1$ layers with 2^s vertices in each layer. For every $i \in \{0, \dots, \ell - 1\}$ each vertex in layer i has 2^v outgoing edges to vertices in layer $i + 1$, one edge for every possible string $x_i \in \{0, 1\}^v$. In addition, layer 0 contains a designated initial vertex, and each vertex in layer ℓ is labeled with 0 or 1. For an (s, v, ℓ) -LBP M and an input $x = (x_1, \dots, x_\ell) \in (\{0, 1\}^v)^\ell$, the computation $M(x)$ is defined by a walk on the graph corresponding to M , starting from the initial vertex in layer 0, and each time advancing to level i along the edge labeled by x_i . The value $M(x)$ is the label of the vertex that is reached in the last layer.

We now define the notion of a pseudorandom generator that ϵ -fools a branching program M . Informally, this means that M can distinguish between the uniform distribution and the output of the generator with probability at most ϵ .

Definition 2.8. A generator $G : \{0, 1\}^m \rightarrow (\{0, 1\}^v)^\ell$ is said to ϵ -fool a layered branching program M if

$$|\Pr[M(G(x)) = 1] - \Pr[M(y) = 1]| \leq \epsilon,$$

where $x \leftarrow \{0, 1\}^m$ and $y \leftarrow (\{0, 1\}^v)^\ell$.

Theorem 2.9 ([23], [18]). For any s, v, ℓ and ϵ there exists a generator $G : \{0, 1\}^m \rightarrow (\{0, 1\}^v)^\ell$ that ϵ -fools any (s, v, ℓ) -LBP, where $m = O(v + \log \ell(s + \log \ell + \log(1/\epsilon)))$, and for any seed $x \in \{0, 1\}^m$ the value $G(x)$ can be computed in time $\text{poly}(s, v, \ell, \log(1/\epsilon))$.

3. A CONSTRUCTION BASED ON GRADUALLY-INCREASING INDEPENDENCE

In this section we present our first family of functions which, as discussed in Section 1.1, is based on replacing $O(\log n / \log \log n)$ -wise independence with “gradually-increasing independence”. The intuition behind the construction is a tree where the root contains all the elements, each level hashes into the next with varying amounts of independence, and the leaves are the desired n bins. More formally, the construction is obtained by concatenating the outputs of $O(\log \log n)$ functions which are *gradually* more independent. Each function $f \in \mathcal{F}$ in our construction consists of d functions h_1, \dots, h_d that are sampled independently, and for any $x \in [u]$ we define $f(x) = h_1(x) \circ \dots \circ h_d(x)$, where we view the output of each h_i as a binary string, and \circ denotes the concatenation operator on binary strings⁴. Going from left to right each function in the above concatenation

³In our setting we are interested in layered branching programs that count the number of balls that are mapped to a specific subset of bins.

⁴Various approaches based on multilevel hashing are fundamental and widely used in the design of data structures (see, for example, [13]).

has a higher level of independence (from $O(1)$ -wise almost independence for h_1 , to $O(\log n / \log \log n)$ -wise almost independence for h_d), and a shorter output length (from $\Omega(\log n)$ bits for h_1 , to $O(\log \log n)$ bits for h_d). This trade-off enables us to represent each of these d functions using $O(\log n)$ bits and to evaluate it in time $O(\log n)$. This construction allows us to prove the following theorem:

Theorem 3.1. For any constant $c > 0$ and integers n and $u = \text{poly}(n)$ there exists a family \mathcal{F} of functions $f : [u] \rightarrow [n]$ satisfying the following properties:

- 1) Each $f \in \mathcal{F}$ is described using $O(\log n \log \log n)$ bits.
- 2) For any $f \in \mathcal{F}$ and $x \in [u]$ the value $f(x)$ can be computed in time $O(\log n \log \log n)$ in the unit cost RAM model.
- 3) There exists a constant $\gamma > 0$ such that for any set $S \subseteq [u]$ of size n it holds that

$$\Pr_{f \leftarrow \mathcal{F}} \left[\max_{i \in [n]} |f^{-1}(i) \cap S| \leq \frac{\gamma \log n}{\log \log n} \right] > 1 - \frac{1}{n^c}.$$

3.1. A Formal Description

To simplify the presentation of our construction we assume that n is a power of two (as otherwise we can choose the number of bins to be the largest power of two which is smaller than n , and this may affect the maximal load by at most a multiplicative factor of two). Let $d = O(\log \log n)$, and for every $i \in [d]$ let \mathcal{H}_i be a family of k_i -wise δ -dependent functions $h_i : [u] \rightarrow \{0, 1\}^{\ell_i}$, where:

- $n_0 = n$, and $n_i = n_{i-1}/2^{\ell_i}$ for every $i \in [d]$.
- $\ell_i = \lfloor (\log n_{i-1})/4 \rfloor$ for every $i \in [d - 1]$, and $\ell_d = \log n - \sum_{i=1}^{d-1} \ell_i$.
- $k_i \ell_i = \Theta(\log n)$ for every $i \in [d - 1]$, and $k_d = \Theta(\log n / \log \log n)$.
- $\delta = \text{poly}(1/n)$.

The exact constants for the construction depend on the error parameter c , and will be fixed by the analysis in Section 3.2. Note that Corollary 2.1 provides such families \mathcal{H}_i where each function $h_i \in \mathcal{H}_i$ is represented using $O(\log u + k_i \ell_i + \log(1/\delta)) = O(\log n)$ bits and evaluated in time $O(\log u + k_i \ell_i + \log(1/\delta)) = O(\log n)$. Each function $f \in \mathcal{F}$ in our construction consists of d functions $h_1 \in \mathcal{H}_1, \dots, h_d \in \mathcal{H}_d$ that are sampled independently and uniformly at random. For any $x \in [u]$ we define $f(x) = h_1(x) \circ \dots \circ h_d(x)$.

3.2. Analyzing the Construction

We view the construction as a tree consisting of $d + 1$ levels that are numbered $0, \dots, d$, where levels $0, \dots, d - 1$ consist of “intermediate” bins, and level d consists of the actual n bins to which the elements are hashed. For a given set $S \subseteq [u]$ of size n , level 0 consists of a single bin containing the n elements of S . Level 1 consists of 2^{ℓ_1} bins, to which the elements of S are hashed using the function h_1 . More generally, each level $i \in \{1, \dots, d - 1\}$ consists of

$2^{\sum_{j=1}^i \ell_j}$ bins, and the elements of each such bin are hashed into $2^{\ell_{i+1}}$ bins in level $i+1$ using the function h_{i+1} .

Recall that we defined $n_0 = n$, and $n_i = n_{i-1}/2^{\ell_i}$ for every $i \in [d]$. The number n_i is the expected number of elements in each bin in level i , and we show that with high probability no bin in levels $0, \dots, d-1$ contains more than $(1+\alpha)^i n_i$ elements, where $\alpha = \Omega(1/\log \log n)$. This will be guaranteed by the following lemma.

Lemma 3.2. *For any $i \in \{0, \dots, d-2\}$, $\alpha = \Omega(1/\log \log n)$, $0 < \alpha_i < 1$, and set $S_i \subseteq [u]$ of size at most $(1+\alpha_i)n_i$ it holds that given $h_{i+1} \leftarrow \mathcal{H}_{i+1}$*

$$\Pr \left[\max_{y \in \{0,1\}^{\ell_{i+1}}} |h_{i+1}^{-1}(y) \cap S_i| \leq (1+\alpha)(1+\alpha_i)n_{i+1} \right] > 1 - \frac{1}{n^{c+1}} .$$

Proof: Fix $y \in \{0,1\}^{\ell_{i+1}}$, let $X = |h_{i+1}^{-1}(y) \cap S_i|$, and assume without loss of generality that $|S_i| \geq \lfloor (1+\alpha_i)n_i \rfloor$. Then X is the sum of $|S_i|$ indicator random variables that are k_{i+1} -wise δ -dependent, and has expectation $\mu = |S_i|/2^{\ell_{i+1}}$. Lemma 2.2 states that

$$\Pr[X > (1+\alpha)\mu] = 2 \left(\frac{2^{2\ell_{i+1}} k_{i+1}}{\alpha^2 |S_i|} \right)^{\frac{k_{i+1}}{2}} + \delta \left(\frac{2^{\ell_{i+1}}}{\alpha} \right)^{k_{i+1}} .$$

We now upper bound each of above two summands separately. For the first one, recall that $\ell_{i+1} \leq (\log n_i)/4$, and combined with the facts that $|S_i| \geq (1+\alpha_i)n_i - 1 \geq n_i$ and $\alpha = \Omega(1/\log \log n)$, this yields

$$2 \left(\frac{2^{2\ell_{i+1}} k_{i+1}}{\alpha^2 |S_i|} \right)^{k_{i+1}/2} \leq \frac{1}{2n^{c+2}} , \quad (3.1)$$

where the last inequality follows from the choice of k_{i+1} and ℓ_{i+1} such that $k_{i+1}\ell_{i+1} = \Omega(\log n)$. This also enables us to upper bound the second summand, noting that for an appropriate choice of $\delta = \text{poly}(1/n)$ it holds that

$$\delta \left(\frac{2^{\ell_{i+1}}}{\alpha} \right)^{k_{i+1}} \leq \frac{1}{2n^{c+2}} . \quad (3.2)$$

Therefore, by combining Equations (3.1) and (3.2), and recalling that $n_{i+1} = n_i/2^{\ell_{i+1}}$ we obtain

$$\Pr[X > (1+\alpha)(1+\alpha_i)n_{i+1}] = \Pr \left[X > (1+\alpha)(1+\alpha_i) \frac{n_i}{2^{\ell_{i+1}}} \right] \leq \frac{1}{n^{c+2}} .$$

The lemma now follows by a union bound over all $y \in \{0,1\}^{\ell_{i+1}}$ (there are at most n such values). ■

We are now ready to prove Theorem 3.1. The description length and evaluation time of our construction were already argued in Section 3.1, and therefore we focus here on the maximal load.

Proof of Theorem 3.1: Fix a set $S \subseteq [u]$ of size n . We begin by inductively arguing that for every level $i \in \{0, \dots, d-1\}$, with probability at least $1 - i/n^{c+1}$ the maximal load in level i is at most $(1+\alpha)^i n_i$ elements per bin, where $\alpha = \Omega(1/\log \log n)$. For $i=0$ this follows by our definition of level 0: it contains a single bin with the $n_0 = n$ elements of S . Assume now that the claim holds for level i and directly apply Lemma 3.2 for each bin in level i with $(1+\alpha_i) = (1+\alpha)^i$. A union bound over all bins in level i implies that with probability at least $1 - (i/n^{c+1} + 1/n^{c+1})$ the maximal load in level $i+1$ is at most $(1+\alpha)^{i+1} n_{i+1}$, and the inductive claim follows. In particular, this guarantees that with probability at least $1 - (d-1)/n^{c+1}$, the maximal load in level $d-1$ is $(1+\alpha)^{d-1} n_{d-1} \leq 2n_{d-1}$, for an appropriate choice of $d = O(\log \log n)$.

Now we would like to upper bound the number n_{d-1} . Note that for every $i \in [d-1]$ it holds that $\ell_i \geq (\log n_{i-1})/4 - 1$, and therefore $n_i = n_{i-1}/2^{\ell_i} \leq 2n_{i-1}^{3/4}$. By simple induction this implies $n_i \leq 2^{\sum_{j=0}^{i-1} (3/4)^j} n^{(3/4)^i} \leq 16n^{(3/4)^i}$. Thus, for an appropriate choice of $d = O(\log \log n)$ it holds that $n_{d-1} \leq \log n$. In addition, the definition of the n_i 's implies that $n_{d-1} = n/2^{\sum_{j=1}^{d-1} \ell_j}$, and therefore $\ell_d = \log n - \sum_{j=1}^{d-1} \ell_j = \log n_{d-1}$.

That is, in level $d-1$ of the tree, with probability at least $1 - (d-1)/n^{c+1}$, each bin contains at most $2n_{d-1} \leq 2 \log n$ elements, and these elements are hashed into n_{d-1} bins using the function h_d . The latter function is k_d -wise δ -dependent, where $k_d = \Omega(\log n / \log \log n)$ and therefore the probability that any $t = \gamma \log n / \log \log n \leq k_d$ elements from level $d-1$ are hashed into any specific bin in level d is at most

$$\binom{2n_{d-1}}{t} \left(\left(\frac{1}{n_{d-1}} \right)^t + \delta \right) \leq \frac{1}{n^{c+3}} ,$$

for an appropriate choice of $t = \gamma \log n / \log \log n$ and $\delta = \text{poly}(1/n)$. This holds for any pair of bins in levels $d-1$ and d , and therefore a union bound over all such bins implies that the probability that there exists a bin in level d with more than t elements is at most $1/n^{c+1}$. This implies that with probability at least $1 - d/n^{c+1} > 1 - 1/n^c$ a randomly chosen function f has a maximal load of $\gamma \log n / \log \log n$. ■

4. A CONSTRUCTION BASED ON GENERATORS FOR SPACE-BOUNDED COMPUTATIONS

Our second construction is based on the observation that any pseudorandom generator which fools small-width branching programs, directly defines a family of functions with the desired maximal load. Specifically, for a universe $[u]$, a generator that produces u blocks of length $\log n$ bits each can be interpreted as a function $f : [u] \rightarrow [n]$, where for any input $x \in [u]$ the value $h(x)$ is defined to be the x th output block of the generator. Fixing a subset $S \subseteq [u]$, the

event in which the load of any particular bin is larger than $t = O(\log n / \log \log n)$ can be recognized by a branching program of width $t + 1 < n$ (the program only needs to count up to t). Assuming the existence of such an explicit generator with seed of length $O(\log n)$ bits implies a family of functions with description length of $O(\log n)$ bits (which is optimal up to a constant factor).

Unfortunately, the best known generator [23] has seed of length $\Omega(\log^2 n)$ bits, which essentially does not give any improvement over using $O(\log n / \log \log n)$ -wise independence. Our construction uses a pseudorandom generator in an inherent way, but instead of generating $O(u \log n)$ bits directly, will only produce $O(2^{\sqrt{\log n}})$ descriptions of $O(\sqrt{\log n})$ -wise independent functions, which we combine into a single function $f : [u] \rightarrow [n]$. The construction is inspired by Lu's generator for combinatorial rectangles [19].

We now describe our family \mathcal{F} . Let \mathcal{H} be a family of k_1 -wise independent functions $h : [u] \rightarrow [\ell]$ for $k_1 = O(1)$ and $\ell = O(2^{\sqrt{\log n}})$, and let \mathcal{G} be a family of k_2 -wise independent functions $g : [u] \rightarrow [n]$ for $k_2 = O(\sqrt{\log n})$. Each function $f \in \mathcal{F}$ consists of a function $h \in \mathcal{H}$ that is sampled uniformly at random, and of ℓ functions $g_1, \dots, g_\ell \in \mathcal{G}$ that are obtained as the output of a pseudorandom generator. The description of each g_j is given by the j th output block of the generator. For any $x \in [u]$ we define $f(x) = g_{h(x)}(x)$.

Using the generator provided by Theorem 2.9, the description length of each f is $O(\log^{3/2} n)$ bits. Moreover, we present a new construction of a pseudorandom generator in which the description of each g_j can be computed in time $O(\sqrt{\log n})$ without increasing the length of the seed. Thus, for any $x \in [u]$ the time required for computing $f(x) = g_{h(x)}(x)$ is $O(\sqrt{\log n})$: the value $h(x)$ can be computed in time $O(k_1) = O(1)$, the description of $g_{h(x)}$ can be computed in time $O(\sqrt{\log n})$, and then the value $g_{h(x)}(x)$ can be computed in time $O(k_2) = O(\sqrt{\log n})$.

Theorem 4.1. *For any constant $c > 0$ and integers n and $u = \text{poly}(n)$ there exists a family \mathcal{F} of functions $f : [u] \rightarrow [n]$ satisfying the following properties:*

- 1) *Each $f \in \mathcal{F}$ is described using $O(\log^{3/2} n)$ bits.*
- 2) *For any $f \in \mathcal{F}$ and $x \in [u]$ the value $f(x)$ can be computed in time $O(\sqrt{\log n})$ in the unit cost RAM model.*
- 3) *There exists a constant $\gamma > 0$ such that for any set $S \subseteq [u]$ of size n ,*

$$\Pr_{f \leftarrow \mathcal{F}} \left[\max_{i \in [n]} |f^{-1}(i) \cap S| \leq \frac{\gamma \log n}{\log \log n} \right] > 1 - \frac{1}{n^c} .$$

The proof of Theorem 4.1 proceeds in three stages. First, in Section 4.1 we analyze a simple family $\widehat{\mathcal{F}}$ where g_1, \dots, g_ℓ are sampled independently and uniformly at random. Then, in Section 4.2 we show that one can replace the descriptions of g_1, \dots, g_ℓ with the output of a pseudorandom

generator with seed length $O(\log^{3/2} n)$ bits. Finally in Section 4.3 we present a new generator that makes it possible to compute the description of each function g_j in time $O(\sqrt{\log n})$ without increasing the length of the seed.

4.1. Analyzing the Basic Construction

We begin by analyzing the family $\widehat{\mathcal{F}}$ in which each function \hat{f} is obtained by sampling $h \in \mathcal{H}$ and $g_1, \dots, g_\ell \in \mathcal{G}$ independently and uniformly at random, and define $\hat{f}(x) = g_{h(x)}(x)$ for any $x \in [u]$. We naturally interpret $\widehat{\mathcal{F}}$ as a two-level process: The function h maps elements into $\ell = O(2^{\sqrt{\log n}})$ first-level bins, and then the elements in each bin j are mapped into n second-level bins using g_j .

When hashing a set $S \subseteq [u]$ of size n , we expect each first-level bin to contain roughly n/ℓ elements, and in Claim 4.2 we observe this holds with high probability. Then, assuming each first-level bin contains roughly n/ℓ elements, Claim 4.3 shows that if the g_j 's are sampled independently and uniformly at random from a family of $O(\sqrt{\log n})$ -wise independent functions, the maximal load in the second-level bins is $O(\log n / \log \log n)$ with high probability.

For a set $S \subseteq [u]$ denote by S_j the subset of S mapped to first level bin j ; i.e. $S_j = S \cap h^{-1}(j)$.

Claim 4.2. *For any set $S \subseteq [u]$ of size n ,*

$$\Pr_{h \leftarrow \mathcal{H}} \left[\max_{j \in [\ell]} |S_j| \leq 2 \cdot \frac{n}{\ell} \right] > 1 - \frac{1}{n^{c+5}} .$$

Proof: For any $j \in [\ell]$ the random variable $|S_j|$ is the sum of n binary random variables that are k_1 -wise independent, and has expectation n/ℓ . Letting $\ell = \beta 2^{\sqrt{\log n}}$ for some constant $\beta > 0$, Lemma 2.2 (for $\delta = 0$ and $t = n/\ell$) guarantees that for an appropriate choice of k_1 ,

$$\begin{aligned} \Pr_{h \leftarrow \mathcal{H}} \left[|S_j| > 2 \cdot \frac{n}{\ell} \right] &\leq 2 \left(\frac{nk_1}{(n/\ell)^2} \right)^{k_1/2} \\ &= 2 \left(\frac{\beta k_1}{2^{\log n - 2\sqrt{\log n}}} \right)^{k_1/2} \leq \frac{1}{n^{c+6}} . \end{aligned}$$

This holds for all $j \in [\ell]$ so a union bound over $\ell \leq n$ yields

$$\Pr_{h \leftarrow \mathcal{H}} \left[\max_{j \in [\ell]} |S_j| > 2 \cdot \frac{n}{\ell} \right] \leq \ell \cdot \frac{1}{n^{c+6}} \leq \frac{1}{n^{c+5}} .$$

■

Claim 4.3. *There exists a constant $\gamma > 0$ such that for any set $S \subseteq [u]$ of size n and for any $i \in [n]$,*

$$\Pr_{\hat{f} \leftarrow \widehat{\mathcal{F}}} \left[\left| \hat{f}^{-1}(i) \cap S \right| \leq \frac{\gamma \log n}{\log \log n} \right] > 1 - \frac{1}{n^{c+4}} .$$

Proof: For any set $S \subseteq [u]$ of size n , Claim 4.2 guarantees that with probability at least $1 - 1/n^{c+5}$, $\max_{j \in [\ell]} |S_j| \leq 2 \cdot \frac{n}{\ell}$. From this point on we condition on the latter event and fix the function h . Let $k_{i,j} = |S_j \cap g_j^{-1}(i)|$ be the number of elements mapped to second level bin i via first level bin j .

The event in which one of the n second-level bins contains more than $t = O(\log n / \log \log n)$ elements is the union of the following two events:

- Event 1: There exists a second-level bin $i \in [n]$ and first level bin j such that $k_{i,j} \geq \alpha \sqrt{\log n}$ for some constant α . If we set g_j to be $\alpha \sqrt{\log n}$ -wise independent, the probability of this event is at most

$$\binom{|S_j|}{\alpha \sqrt{\log n}} \cdot \left(\frac{1}{n}\right)^{\alpha \sqrt{\log n}} \leq \binom{|S_j|}{n}^{\alpha \sqrt{\log n}} \leq \frac{1}{n^{c+7}}$$

for an appropriate choice of $\ell = O(2^{\sqrt{\log n}})$ and α . There are $n\ell < n^2$ pairs of bins, thus a union bound implies this occurs with probability at most $1/n^{c+5}$.

- Event 2: Some second-level bin $i \in [n]$ has $t = O(\log n / \log \log n)$ elements with $k_{i,j} \leq \alpha \sqrt{\log n}$ for all j . Since g_1, \dots, g_ℓ are sampled independently from a family that is $\alpha \sqrt{\log n}$ -wise independent, the probability of this event is at most

$$\binom{n}{t} \left(\frac{1}{n}\right)^t \leq \left(\frac{ne}{tn}\right)^t \leq \frac{1}{n^{c+6}}$$

for an appropriate choice of $t = O(\log n / \log \log n)$. This holds for any second-level bin $i \in [n]$, thus a union bound over all n bins implies this event occurs with probability at most $1/n^{c+5}$.

Combining all of the above, we obtain there exists a constant γ such that for all sufficiently large n ,

$$\Pr_{\hat{f} \leftarrow \hat{\mathcal{F}}} \left[\max_{i \in [n]} |\hat{f}^{-1}(i) \cap S| \leq \frac{\gamma \log n}{\log \log n} \right] > 1 - \frac{1}{n^{c+4}} .$$

4.2. Derandomizing the Basic Construction

As discussed above, the key observation that allows the derandomization $g_1, \dots, g_\ell \in \mathcal{G}$ is the fact that the event in which the load of any particular bin is larger than $t = O(\log n / \log \log n)$ can be recognized in $O(\log n)$ space (and, more accurately, in $O(\log t)$ space). Specifically, fix a set $S \subseteq [u]$ of size n , a second-level bin $i \in [n]$, and a function $h \in \mathcal{H}$. Consider the layered branching program $M_{S,h,i}$ that has $\ell + 1$ layers each of which contains roughly n vertices, where every layer $j \in [\ell]$ takes as input the description of the function g_j and keeps count of the number of elements from S that are mapped to bin i using the functions h, g_1, \dots, g_j . In other words, the j th layer adds to the count the number of elements $x \in S$ such that $h(x) = j$ and $g_j(x) = i$. Each vertex in the final layer is labeled with 0 and 1 depending on whether the count has exceeded t (note there is no reason to keep on counting beyond t , thus it suffices to have only t vertices in each layer).

Let $G : \{0, 1\}^m \rightarrow (\{0, 1\}^v)^\ell$ be a generator that ϵ -fools any (s, v, ℓ) -LBP, where $\epsilon = 1/n^{c+4}$, $s = O(\log n)$, $v = O(k_2 \log n) = O(\log^{3/2} n)$, and $\ell = O(2^{\sqrt{\log n}})$.

Theorem 2.9 provides an explicit construction of such a generator with a seed of length $m = O(v + \log \ell \cdot (s + \log \ell + \log(1/\epsilon))) = O(\log^{3/2} n)$. For any seed $x \in \{0, 1\}^m$ we use $G(x) = (g_1, \dots, g_\ell) \in (\{0, 1\}^v)^\ell$ for providing the descriptions of the function g_1, \dots, g_ℓ .

By combining Claim 4.3 with the fact that G is a generator that ϵ -fools any (s, v, ℓ) -LBP we obtain the following claim:

Claim 4.4. *There exists a constant $\gamma > 0$ such that for any set $S \subseteq [u]$ of size n ,*

$$\Pr_{f \leftarrow \mathcal{F}} \left[\max_{i \in [n]} |f^{-1}(i) \cap S| \leq \frac{\gamma \log n}{\log \log n} \right] > 1 - \frac{1}{n^c} .$$

Proof: Let $\gamma > 0$ be the constant specified by Claim 4.3. Then for any set $S \subseteq [u]$ of size n and $i \in [n]$,

$$\begin{aligned} & \Pr_{f \leftarrow \mathcal{F}} \left[|f^{-1}(i) \cap S| > \frac{\gamma \log n}{\log \log n} \right] \\ & \leq \Pr_{h \leftarrow \mathcal{H}} \left[\Pr_{g_1, \dots, g_\ell \leftarrow \mathcal{G}} [M_{S,h,i}(g_1, \dots, g_\ell) = 1] \right] + \frac{1}{n^{c+4}} \\ & < \frac{2}{n^{c+4}} . \end{aligned}$$

The first inequality follows since G is a generator that $1/n^{c+4}$ -fools $M_{S,h,i}$, and the second follows from Claim 4.3. A union bound over all bins $i \in [n]$ yields

$$\Pr_{f \leftarrow \mathcal{F}} \left[\max_{i \in [n]} |f^{-1}(i) \cap S| \geq \frac{c \log n}{\log \log n} \right] \leq \frac{1}{n^c} .$$

4.3. A More Efficient Generator

As shown in Section 4.2, we can instantiate our construction with any generator $G : \{0, 1\}^m \rightarrow (\{0, 1\}^v)^\ell$ that ϵ -fools any (s, v, ℓ) -LBP, where $s = O(\log n)$, $v = O(\log^{3/2} n)$, $\ell = O(2^{\sqrt{\log n}})$, and $\epsilon = 1/n^{c+4}$. The generator constructed by Impagliazzo et. al. [18] (following [23]), whose parameters we stated in Theorem 2.9, provides one instantiation, but the time it requires to compute each v -bit output block seems at least logarithmic. Here we construct a more efficient generator for our parameters, where each v -bit output block can be computed in time $O(\sqrt{\log n})$ without increasing the length of the seed.

The generator we propose uses as a building blocks the generators constructed by Nisan [23] and by Nisan and Zuckerman [24]. We first provide a high-level description of these two generators before describing our generator.

4.3.1. Nisan's Generator: Let \mathcal{H} be a family of pairwise independent functions $h : \{0, 1\}^{v_2} \rightarrow \{0, 1\}^{v_2}$. For every integer $k \geq 0$ Nisan [23] constructed a generator $G_N^{(k)} : \{0, 1\}^{v_2} \times \mathcal{H}^k \rightarrow (\{0, 1\}^{v_2})^{2^k}$ that is defined recursively by $G_N^{(0)}(x) = x$, and $G_N^{(k)}(x, h_1, \dots, h_k) = G_N^{(k-1)}(x, h_1, \dots, h_{k-1}) \circ G_N^{(k-1)}(h_k(x), h_1, \dots, h_k)$, where \circ denotes the concatenation operator. For any integers v_2 and

k , Nisan proved that $G_N^{(k)}$ is a generator that 2^{-v_2} -fools any $(v_2, v_2, 2^k)$ -LBP.

When viewing the output of the generator as the concatenation of 2^k blocks of length v_2 bits each, we observe that each block can be computed by evaluating k pairwise independent functions. In our setting we are interested in $v_2 = O(\log n)$ and $2^k = O(2^{\sqrt{\log n}})$, and in this case each output block can be computed in time $O(\sqrt{\log n})$. Formally, from Nisan's generator we obtain the following corollary:

Corollary 4.5. *For any $s_2 = O(\log n)$, $v_2 = O(\log n)$, $\ell_2 = O(2^{\sqrt{\log n}})$, and $\epsilon = \text{poly}(1/n)$, there exists a generator $G_N : \{0, 1\}^{m_2} \rightarrow (\{0, 1\}^{v_2})^{\ell_2}$ that ϵ -fools any (s_2, v_2, ℓ_2) -LBP, where $m_2 = O(\log^{3/2} n)$. In the unit cost RAM model with a word size of $w = \Omega(\log n)$ bits each v_2 -bit output block can be computed in time $O(\sqrt{\log n})$.*

4.3.2. The Nisan-Zuckerman Generator and an Efficient Instantiation: Given a (k, ϵ) -extractor $\text{Ext} : \{0, 1\}^{t_1} \times \{0, 1\}^{d_1} \rightarrow \{0, 1\}^{v_1}$ Nisan and Zuckerman [24] constructed a generator $G_{\text{NZ}}^{\text{Ext}} : \{0, 1\}^{t_1} \times (\{0, 1\}^{d_1})^{\ell_1} \rightarrow (\{0, 1\}^{v_1})^{\ell_1}$ that is defined as $G_{\text{NZ}}^{\text{Ext}}(x, y_1, \dots, y_{\ell_1}) = \text{Ext}(x, y_1) \circ \dots \circ \text{Ext}(x, y_{\ell_1})$, where \circ denotes the concatenation operator. When viewing the output of the generator as the concatenation of ℓ_1 blocks of length v_1 bits each, we observe that the time to compute each block is the time to compute the extractor Ext . In our setting we are interested in $t_1 = O(\log^{3/2} n)$, $v_1 = O(\log^{3/2} n)$, $k = t_1 - O(\log n)$, and $\epsilon = \text{poly}(1/n)$, and in Lemma 4.7 we construct an extractor that has a seed of length $d_1 = O(\log n)$ bits and can be computed in time $O(\sqrt{\log n})$. As a corollary, from the Nisan-Zuckerman generator when instantiated with our extractor we obtain:

Corollary 4.6. *For any $s_1 = O(\log n)$, $v_1 = O(\log^{3/2} n)$, $\ell_1 = O(2^{\sqrt{\log n}})$, and $\epsilon = \text{poly}(1/n)$, there exists a generator $G_{\text{NZ}}^{\text{Ext}} : \{0, 1\}^{m_1} \rightarrow (\{0, 1\}^{v_1})^{\ell_1}$ that ϵ -fools any (s_1, v_1, ℓ_1) -LBP, where $m_1 = O(\log^{3/2} n + 2^{\sqrt{\log n}} \cdot \log n)$. Moreover, there exists an extractor Ext such that in the unit cost RAM model with a word size of $w = \Omega(\log n)$ bits each v_1 -bit output block of the generator $G_{\text{NZ}}^{\text{Ext}}$ can be computed in time $O(\sqrt{\log n})$.*

The following lemma presents the extractor that we use for instantiating the Nisan-Zuckerman generator.

Lemma 4.7. *Let $t_1 = \Theta(\log^{3/2} n)$, $\Delta = O(\log n)$ and $\epsilon = \text{poly}(1/n)$. There exists a $(t_1 - \Delta, \epsilon)$ -extractor $\text{Ext} : \{0, 1\}^{t_1} \times \{0, 1\}^{d_1} \rightarrow \{0, 1\}^{v_1}$, where $d_1 = O(\log n)$ and $v_1 = \Omega(\log^{3/2} n)$, that can be computed in time $O(\sqrt{\log n})$ in the unit cost RAM model with a word size of $w = \Omega(\log n)$ bits.*

Proof: Given a random variable X over $\{0, 1\}^{t_1}$ we partition it into $T = t_1/z$ consecutive blocks $X = X_1 \circ \dots \circ X_T$ each of length z bits, where $z = \lceil 2(\Delta + 3 \log(2T/\epsilon)) \rceil = O(\log n)$. Without loss of generality we assume that z divides t_1 , and otherwise we ignore the last block. Corollary 2.5 guarantees that X is a $(T, z - \Delta - \log(2T/\epsilon), \epsilon/2T)$ -block source. Let \mathcal{H} be a family of pair-wise independent functions $h : \{0, 1\}^z \rightarrow \{0, 1\}^{z'}$, where $z' = \lfloor z - \Delta - 3 \log(2T/\epsilon) \rfloor = \Omega(\log n)$ and each $h \in \mathcal{H}$ is described by $d_1 = O(\log n)$ bits. We define an extractor $\text{Ext} : \{0, 1\}^{t_1} \times \mathcal{H} \rightarrow \{0, 1\}^{Tz'}$ by applying a randomly chosen $h \in \mathcal{H}$ to each of the T blocks of the source. That is,

$$\text{Ext}(x_1 \circ \dots \circ x_T, h) = h(x_1) \circ \dots \circ h(x_T) .$$

Lemma 2.7 implies the distribution $(h, h(x_1), \dots, h(x_T))$ is ϵ -close to the distribution (h, y_1, \dots, y_T) , where $h \leftarrow \mathcal{H}$, $(x_1, \dots, x_T) \leftarrow X$, and $(y_1, \dots, y_T) \leftarrow (\{0, 1\}^{z'})^T$. In addition, in the unit cost RAM model with a word size of $w = \Omega(\log n)$ bits each application of h can be done in constant time, thus the extractor can be computed in time $T = O(\sqrt{\log n})$. Finally, note that the number of outputs bits is $Tz' = t_1 z'/z = \Omega(\log^{3/2} n)$. ■

4.3.3. Our Generator: Recall that we are interested in a generator $G : \{0, 1\}^m \rightarrow (\{0, 1\}^v)^\ell$ that ϵ -fools any (s, v, ℓ) -LBP, where $s = O(\log n)$, $v = O(\log^{3/2} n)$, $\ell = O(2^{\sqrt{\log n}})$, and $\epsilon = 1/n^{c+4}$. Let $G_{\text{NZ}} : \{0, 1\}^{m_1} \rightarrow (\{0, 1\}^v)^\ell$ be the Nisan-Zuckerman generator that is given by Corollary 4.6 that $\epsilon/2$ -fools any (s, v, ℓ) -LBP, where $m_1 = O(\log^{3/2} n + 2^{\sqrt{\log n}} \cdot \log n)$. In addition, let $G_N : \{0, 1\}^{m_2} \rightarrow (\{0, 1\}^{v_2})^\ell$ be Nisan's generator that is given by Corollary 4.5 that $\epsilon/2$ -fools any (s, v_2, ℓ) -LBP, where $v_2 = O(\log n)$ and $m_2 = O(\log^{3/2} n)$. We define a generator G as follows: $G(x_1, x_2) = G_{\text{NZ}}(x_1, G_N(x_2))$.

That is, given a seed (x_1, x_2) it first computes the output (y_1, \dots, y_ℓ) of Nisan's generator using the seed x_2 , and then it computes the output $\text{Ext}(x_1, y_1) \circ \dots \circ \text{Ext}(x_1, y_{\ell_1})$ of the Nisan-Zuckerman generator. Observe that the time to compute the i th v -bit output block is the time to compute the i th output block for both generators, which is $O(\sqrt{\log n})$. In addition, note that the length of seed is $O(\log^{3/2} n)$ bits since each of x_1 and x_2 is of length $O(\log^{3/2} n)$ bits. Thus, it only remains to prove that G indeed ϵ -fools any (s, v, ℓ) -LBP. This follows from the construction, and the proof is contained in the full version of this paper [6].

Lemma 4.8. *For the parameters s, v, ℓ, m , and ϵ specified above, G is a generator that ϵ -fools any (s, v, ℓ) -LBP.*

5. EXTENSIONS AND OPEN PROBLEMS

Applications: Our constructions can be employed for storing elements using linear probing, guaranteeing an inser-

tion time of $O(\log n)$ with high probability. An interesting problem is whether these, or similar techniques, can be used to construct small hash families that are suitable for other applications. For instance, $O(\log n)$ -wise independence is known to suffice for two-choice hashing [3], [31], cuckoo hashing [27], and more. Existing constructions have focused on simplicity and fast computation [11], [32], [25], [28], albeit with a significant increase to the description length.

Augmenting our constructions with k -wise independence: Our constructions can be augmented to offer $O(\log \log n)$ -wise and $O(\log^{1/2} n)$ -wise independence respectively without affecting their description length and evaluation time. Specifically, a function f resulting from either the first or second construction can be modified to $f(x) + h(x) \bmod n$, where h is sampled from a family of $O(\log \log n)$ - or $O(\log^{1/2} n)$ -wise independent functions respectively. Our analysis easily extends to this case, and the resulting functions enjoy the level of independence offered by h . By implementing h appropriately, the description length and evaluation time of our constructions are not affected. This may be useful, for example, in any application that involves tail bounds for limited independence.

A time-space lower bound: Our constructions offer two different trade-offs between the description length and the evaluation time. It would be interesting to prove a lower bound on the time-space trade-off of any family that guarantees a maximal load of $O(\log n / \log \log n)$ with high probability when hashing n balls into n bins. For example, can we rule out constructions that are optimal in both description length and evaluation time?

ACKNOWLEDGMENT

We thank Moni Naor for many inspiring discussions.

REFERENCES

- [1] N. Alon, M. Dietzfelbinger, P. Miltersen, E. Petrank, and G. Tardos, "Linear hash functions," *J. of the ACM*, vol. 46.
- [2] N. Alon, O. Goldreich, J. Håstad, and R. Peralta, "Simple construction of almost k -wise independent random variables," *RANDOM*, vol. 3.
- [3] Y. Azar, A. Z. Broder, A. R. Karlin, and E. Upfal, "Balanced allocations," *SIAM J. on Computing*, vol. 29.
- [4] M. Bellare and J. Rompel, "Randomness-efficient oblivious sampling," in *FOCS*, 1994, pp. 276–287.
- [5] L. Carter and M. N. Wegman, "Universal classes of hash functions," *J. of Computer and System Sciences*, vol. 18, pp. 143–154, 1979.
- [6] L. E. Celis, O. Reingold, G. Segev, and U. Wieder, "Balls and bins: Smaller hash families and faster evaluation (full version)," ECCC report TR11-068, 2011.
- [7] B. Chor and O. Goldreich, "Unbiased bits from sources of weak randomness and probabilistic communication complexity," *SIAM J. on Computing*, vol. 17.
- [8] M. Dietzfelbinger and F. Meyer auf der Heide, "A new universal class of hash functions and dynamic hashing in real time," in *ICALP*, 1990, pp. 6–19.
- [9] M. Dietzfelbinger and R. Pagh, "Succinct data structures for retrieval and approximate membership," in *ICALP*, 2008, pp. 385–396.
- [10] M. Dietzfelbinger and M. Rink, "Applications of a splitting trick," in *ICALP*, 2009, pp. 354–365.
- [11] M. Dietzfelbinger and P. Woelfel, "Almost random graphs with simple hash functions," in *STOC*, 2003, pp. 629–638.
- [12] D. P. Dubhashi and A. Panconesi, *Concentration of Measure for the Analysis of Randomized Algorithms*. Cambridge University Press, 2009.
- [13] M. L. Fredman, J. Komlós, and E. Szemerédi, "Storing a sparse table with $O(1)$ worst case access time," *J. of the ACM*, vol. 31, pp. 538–544, 1984.
- [14] O. Goldreich and A. Wigderson, "Tiny families of functions with random properties: A quality-size trade-off for hashing," *RANDOM*, vol. 11.
- [15] T. Hagerup, "Sorting and searching on the word RAM," in *STOC*, 1998, pp. 366–398.
- [16] T. Hagerup, P. B. Miltersen, and R. Pagh, "Deterministic dictionaries," *J. of Algorithms*, vol. 41, pp. 69–85, 2001.
- [17] R. Impagliazzo, L. A. Levin, and M. Luby, "Pseudo-random generation from one-way functions," in *STOC*, 1989, pp. 12–24.
- [18] R. Impagliazzo, N. Nisan, and A. Wigderson, "Pseudorandomness for network algorithms," in *STOC*, 1994, pp. 356–364.
- [19] C.-J. Lu, "Improved pseudorandom generators for combinatorial rectangles," *Combinatorica*, pp. 417–434, 2002.
- [20] P. B. Miltersen, "Cell probe complexity - a survey," in *FSTTCS*, 1999.
- [21] M. Mitzenmacher and S. Vadhan, "Why simple hash functions work: Exploiting the entropy in a data stream," in *SODA*, 2008, pp. 746–755.
- [22] J. Naor and M. Naor, "Small-bias probability spaces: Efficient constructions and applications," *SIAM J. on Computing*, vol. 22.
- [23] N. Nisan, "Pseudorandom generators for space-bounded computation," *Combinatorica*, vol. 12.
- [24] N. Nisan and D. Zuckerman, "Randomness is linear in space," *J. of Computer and System Sciences*, vol. 52.
- [25] A. Pagh and R. Pagh, "Uniform hashing in constant time and optimal space," *SIAM J. on Computing*, vol. 38.
- [26] A. Pagh, R. Pagh, and M. Ružić, "Linear probing with constant independence," in *STOC*, 2007, pp. 318–327.
- [27] R. Pagh and F. F. Rodler, "Cuckoo hashing," *J. of Algorithms*, vol. 51, pp. 122–144, 2004.
- [28] M. Pătraşcu and M. Thorup, "The power of simple tabulation hashing," *STOC*, 2011.
- [29] J. P. Schmidt and A. Siegel, "The analysis of closed hashing under limited randomness," in *STOC*, 1990, pp. 224–234.
- [30] A. Siegel, "On universal classes of extremely random constant-time hash functions," *SIAM J. on Computing*, vol. 33.
- [31] B. Vöcking, "How asymmetry helps load balancing," *J. of the ACM*, vol. 50.
- [32] P. Woelfel, "Asymmetric balanced allocation with simple hash functions," in *SODA*, 2006, pp. 424–433.
- [33] D. Zuckerman, "Simulating BPP using a general weak random source," *Algorithmica*, vol. 16.