

# Min-Wise Independent Permutations

(Extended abstract)

Andrei Z. Broder\*

Moses Charikar†

Alan M. Frieze‡

Michael Mitzenmacher§

## Abstract

We define and study the notion of min-wise independent families of permutations. We say that  $\mathcal{F} \subseteq S_n$  is *min-wise independent* if for any set  $X \subseteq [n]$  and any  $x \in X$ , when  $\pi$  is chosen at random in  $\mathcal{F}$  we have

$$\Pr(\min\{\pi(X)\} = \pi(x)) = \frac{1}{|X|}.$$

In other words we require that all the elements of any fixed set  $X$  have an equal chance to become the minimum element of the image of  $X$  under  $\pi$ .

Our research was motivated by the fact that such a family (under some relaxations) is essential to the algorithm used in practice by the AltaVista web index software to detect and filter near-duplicate documents. However, in the course of

our investigation we have discovered interesting and challenging theoretical questions related to this concept – we present the solution to some of them and we list the rest as open problems.

## 1 Introduction

The classic analysis of hashing schemes often entails the assumption that the hash functions used are random. More precisely, the assumption is that keys belonging to a universe  $\mathcal{U}$  are hashed into a table of size  $M$  by choosing a function  $h$  uniformly at random among all the functions  $\mathcal{U} \rightarrow [M]$ . (The notation  $[M]$  stands for the set  $\{0, \dots, M-1\}$ . This is slightly non-standard, but convenient for our purposes.) This assumption is impractical since just specifying such a function requires  $|\mathcal{U}| \log(M)$  bits, which usually far exceeds the available storage.

Fortunately in most cases heuristic hash functions behave very closely to the expected behavior of random hash functions; but there are cases when rigorous probabilistic guarantees are necessary. For instance, various adaptive hashing schemes presume that a hash function with certain prescribed properties can be found in constant expected time. This holds if the function is chosen uniformly at random from all possible functions until a suitable one is found but not necessarily if the search is limited to a smaller set of functions.

This situation has led Carter and Wegman [8] to the concept of *universal hashing*. A family of hash functions  $\mathcal{H}$  is called *weakly universal* if for any pair of distinct elements  $x_1, x_2 \in \mathcal{U}$ , if  $h$  is chosen uniformly at random from  $\mathcal{H}$  then

$$\Pr(h(x_1) = h(x_2)) \leq \frac{1}{|M|} \quad (1)$$

and is called (*strongly*) *universal* or *pair-wise independent* if for any pair of distinct elements  $x_1, x_2 \in \mathcal{U}$  and arbitrary  $y_1, y_2 \in [M]$

$$\Pr(h(x_1) = y_1 \text{ and } h(x_2) = y_2) = \frac{1}{|M|^2}. \quad (2)$$

It turns out that in many situations the analysis of various hashing schemes can be completed under the weaker assumption that  $h$  is chosen uniformly at random from a universal

\*Digital SRC, 130 Lytton Avenue, Palo Alto, CA 94301, USA. E-mail: broder@pa.dec.com.

†Computer Science Department, Stanford University, CA 94305, USA. E-mail: mooses@cs.stanford.edu. Part of this work was done while this author was a summer intern at Digital SRC. Supported by the Pierre and Christine Lamond Fellowship and in part by an ARO MURI Grant DAAH04-96-1-0007 and NSF Award CCR-9357849, with matching funds from IBM, Schlumberger Foundation, Shell Foundation, and Xerox Corporation.

‡Department of Mathematical Sciences, Carnegie Mellon University, Pittsburgh, Pennsylvania 15213, USA. Part of this work was done while this author was visiting Digital SRC. Supported in part by NSF grant CCR9530974. E-mail: aflp@andrew.cmu.edu

§Digital SRC, 130 Lytton Avenue, Palo Alto, CA 94301, USA. E-mail: michaelm@pa.dec.com.

family, rather than the assumption that  $h$  is chosen uniformly at random from among all possible function. In other words limited randomness suffices. Furthermore there exist universal families of size  $O(|M|^2)$  that can be easily implemented in practice. Thus, universal hash functions are very useful in the design of adaptive hash schemes (see e.g. [7, 9]) and are actually used in commercial high-performance products (see e.g. [13]). Not only that, but the concept of pairwise independence has important theoretical application. (See the excellent survey by Luby and Wigderson [11].)

It is often convenient to consider permutations rather than functions. Let  $S_n$  be the set of all permutations of  $[n]$ . We say that a family of permutations  $\mathcal{F} \subseteq S_n$  is *pair-wise independent* if for any  $\{x_1, x_2, y_1, y_2\} \subseteq [n]$  with  $x_1 \neq x_2$  and  $y_1 \neq y_2$

$$\Pr(\pi(x_1) = y_1 \text{ and } \pi(x_2) = y_2) = \frac{1}{n(n-1)} \quad (3)$$

In a similar vein, in this paper, we say that  $\mathcal{F} \subseteq S_n$  is *exactly min-wise independent* (or just *min-wise independent* where the meaning is clear) if for any set  $X \subseteq [n]$  and any  $x \in X$ , when  $\pi$  is chosen at random in  $\mathcal{F}$  we have

$$\Pr(\min\{\pi(X)\} = \pi(x)) = \frac{1}{|X|}. \quad (4)$$

In other words we require that all the elements of any fixed set  $X$  have an equal chance to become the minimum element of the image of  $X$  under  $\pi$ . Unless otherwise stated we shall assume that  $\pi$  is chosen uniformly at random in  $\mathcal{F}$ ; otherwise, we shall say  $\pi$  is chosen with a *biased* distribution  $\mu$ . Uniform distributions are natural in this setting, since in practice they are simple to represent.

As explained below, this definition is motivated by the fact that such a family (under some relaxations) is essential to the algorithm currently used in practice by the AltaVista Web indexing software [12] to detect and filter near-duplicate documents.

The Web [4] has undergone exponential growth since its birth, and this has led to the proliferation of documents that are identical or near identical. Experiments indicate that over 20% of the publicly available documents on the web are duplicates or near-duplicates. These documents arise innocently (e.g. local copies of popular documents, mirroring), maliciously (e.g., “spammers” and “robot traps”), and erroneously (spider mistakes). In any case they represent a serious problem for indexing software for two main reasons: first, indexing of duplicates wastes expensive resources and second, users are seldom interested in seeing documents that are “roughly the same” in response to their queries.

This informal concept does not seem to be well captured by any of the standard distances defined on strings (Hamming, Levenshtein, etc.). Furthermore the computation of these distances usually requires the pairwise comparison of entire documents. For a very large collection of documents this is not feasible, and a sampling mechanism per document is necessary.

It turns out that the problem can be reduced to a set intersection problem by a process called *shingling*. (See [6, 5] for details.) Via shingling each document  $D$  gets an associated set  $S_D$ . For the purpose of the discussion here we can view  $S_D$  as a set of natural numbers. (The size of  $S_D$  is about equal to the number of words in  $D$ .) The *resemblance*  $r(A, B)$  of two documents,  $A$  and  $B$ , is defined as

$$r(A, B) = \frac{|S_A \cap S_B|}{|S_A \cup S_B|}.$$

Experiments seem to indicate that high resemblance (that is, close to 1) captures well the informal notion of “near-duplicate” or “roughly the same”.

To compute the resemblance of two documents it suffices to keep for each document a relatively small, fixed size *sketch*. The sketches can be computed fairly fast (linear in the size of the documents) and given two sketches the resemblance of the corresponding documents can be computed in linear time in the size of the sketches.

This is done as follows. Assume that for all documents of interest  $S_D \subseteq \{1, \dots, n\}$ . (In practice  $n = 2^{64}$ .) Let  $\pi$  be chosen uniformly at random over  $S_n$ , the set of permutations of  $[n]$ . Then

$$\begin{aligned} \Pr(\min\{\pi(S_A)\} = \min\{\pi(S_B)\}) \\ = \frac{|S_A \cap S_B|}{|S_A \cup S_B|} = r(A, B). \end{aligned} \quad (5)$$

Hence, we can choose say, 100 independent random permutations  $\pi_1, \dots, \pi_{100}$ . For each document  $D$ , we store the list

$$\bar{S}_A = (\min\{\pi_1(S_A)\}, \min\{\pi_2(S_A)\}, \dots, \min\{\pi_{100}(S_A)\}).$$

Then we can readily estimate the resemblance of  $A$  and  $B$  by computing how many corresponding elements in  $\bar{S}_A$  and  $\bar{S}_B$  are common.

In practice, as in the case of hashing discussed above, we have to deal with the sad reality that it is impossible to choose  $\pi$  uniformly at random in  $S_n$ . We are thus led to consider smaller families of permutations that still satisfy the min-wise independence condition given by equation (4), since min-wise independence is necessary and sufficient for equation (5) to hold.

In practice we can allow certain relaxations. First we can accept small relative errors. We say that  $\mathcal{F} \subseteq S_n$  is *approximately min-wise independent* if for any set  $X \subseteq [n]$  and any  $x \in X$ , when  $\pi$  is chosen at random in  $\mathcal{F}$  we have

$$\left| \Pr(\min\{\pi(X)\} = \pi(x)) - \frac{1}{|X|} \right| \leq \frac{\epsilon}{|X|}. \quad (6)$$

In other words we require that all the elements of any fixed set  $X$  have only an almost equal chance to become the minimum element of the image of  $X$  under  $\pi$ . The expected relative error made in evaluating resemblance using approximately min-wise independent families is less than  $\epsilon$ .

Second the sets of interest are usually much smaller than  $n$ . (For the situation discussed above the typical set has size 1000 while  $n = 2^{64}$ .) We say that  $\mathcal{F} \subseteq S_n$  is *restricted min-wise independent* if for any set  $X \subseteq [n]$  with  $|X| \leq k$  and any  $x \in X$ , when  $\pi$  is chosen at random in  $\mathcal{F}$  we have

$$\Pr(\min\{\pi(X)\} = \pi(x)) = \frac{1}{|X|}, \quad |X| \leq k. \quad (7)$$

Thirdly and finally, it turns out that whether the distribution on the family  $\mathcal{F}$  is uniform or not leads to qualitatively different results.

Ultimately we are interested in practical families of permutations. Hence we first study what is the minimum size of a family that satisfies various combinations of requirements. Clearly if the minimum size is exponential no practical solution exists. It turns out that the exact min-wise property generally necessitates exponential size but that the approximate property can be satisfied by polynomial size families. The complete synopsis of our results is given in Table 1. The entries for which we have no bounds beyond those implied by other entries in the table are marked “?” and the entries for which we have no non-trivial bounds are marked “???”.

Starting from the opposite end we study how good is the performance provided by various families that are easily implementable in software. We consider pair-wise independent families, for which there are numerous practical implementations. In particular we are interested in linear transformations, since they are used in the AltaVista implementation and are known to perform better in some situations than other pair-wise independent families (see [1]).

The way we evaluate this performance is to consider a set  $X$  and study the distribution of the minimum of the image of  $X$ . It suffices to examine the two elements that are respectively most likely and least likely to become the minimum since all the other elements will become the minimum with a probability in between the extremal values. We consider two situations: when  $X$  is chosen to be the worst set (farthest from uniform) with regard to the property of interest; and when  $X$  is chosen uniformly at random, in which case we look for the expected value of the bound over the random choices of  $X$ . The synopsis of our answers is given in Table 2, where we follow the same convention as before regarding the use of “?” and “???”.

## 2 Exact Min-Wise Independence

In this section, we provide bounds for the size of families that are exactly min-wise independent. We begin by determining a lower bound, demonstrating that the size of the family  $\mathcal{F}$  must grow exponentially with  $n$ .

**Theorem 1** *Let  $\mathcal{F}$  be min-wise independent. Then  $|\mathcal{F}|$  is at least as large as the least common multiple (lcm) of the numbers  $1, 2, \dots, n$ , and hence  $|\mathcal{F}| \geq e^{n-o(n)}$ .*

*Proof:* Let  $X$  be a subset of  $[n]$  with  $|X| = j$ . Each element of  $X$  must be the minimum under the family  $\mathcal{F}$  the same

number of times, so  $j$  must divide  $|\mathcal{F}|$ . This holds for every  $j \in \{1, 2, \dots, n\}$ , so the lcm of  $\{1, 2, \dots, n\}$  must divide  $|\mathcal{F}|$ . That the lcm of the first  $n$  numbers has size  $e^{n-o(n)}$  is a well known fact of number theory [3, p. 76].  $\square$

**Remark 1** *This proof also gives a lower bound of  $e^{k-o(k)}$  for restricted min-wise independent families. Also, note that the proof does not require that the members of  $\mathcal{F}$  be distinct. Hence the theorem holds even if  $\mathcal{F}$  contains duplicates of some permutations.*

We now describe a min-wise independent family of size less than  $4^n$ , which is significantly smaller than the trivial bound of  $n!$  and of the same form as the lower bound given above.

**Theorem 2** *There exists a min-wise independent family  $\mathcal{F}$  of size less than  $4^n$ .*

*Proof:* We initially assume for convenience, that  $n = 2^r$  for some  $r$ . We construct the family of permutations recursively in stages. In the first stage, we divide the set  $[n]$  into two equal halves, the top and the bottom. At the first stage, there are  $\binom{n}{n/2}$  ways to partition the set. Each of these can be described by an  $n$  bit string with exactly  $n/2$  ones in it. Element  $i$  goes in the top half if and only if the bit string has a 1 in the  $i$ th position. We proceed to partition each half. Again this can be done by choosing a  $n/2$  bit string with  $n/4$  ones in it. There are  $\binom{n/2}{n/4}$  such strings. Importantly, we use the same string for each half. At the  $i$ th stage, we have the set divided into  $2^{i-1}$  parts each of size  $n/2^{i-1}$ . We partition each part into two halves by choosing a  $n/2^{i-1}$  bit string with  $n/2^i$  ones and using this string to define the partition for each of the  $2^{i-1}$  parts. We continue in this way until each part has size 1. This process produces a permutation of the set in a natural way, with the topmost element receiving the smallest number in the permutation.

The property that each element is the minimum with the correct probability can be verified directly by calculation. More intuitively, when we split  $[n]$  into two halves, every element of  $X$  has an equal chance to go to the upper half or to the lower half; furthermore, all elements of  $X$  now in the top half are equally likely to eventually become the topmost element of  $X$  (by induction). If no elements of  $X$  are in the top half, then all lie in the bottom, and again (by induction) all are equally likely to become eventually the topmost.

The number of permutations in this family is

$$\prod_{i=1}^{\log n} \binom{n/2^{i-1}}{n/2^i}.$$

A simple calculation shows that the size of this family is  $4^{n-O(\log^2 n)}$ .

We can easily remove the assumption that  $n$  is a power of 2. We leave this as an exercise to the diligent reader.  $\square$

Family type	Upper bound	Lower bound
Exactly min-wise, uniform distrib on $\mathcal{F}$	$4^n$	$e^{n-o(n)}$
Exactly min-wise, biased distrib on $\mathcal{F}$	$n2^{n-1} - 1$	$\Omega(\sqrt{n}2^n)$
Exactly min-wise, restricted, uniform distrib on $\mathcal{F}$	?	$e^{k-o(k)}$
Exactly min-wise, restricted, biased distrib on $\mathcal{F}$	$\sum_{j < k} j \binom{n}{j}$	$\Omega\left(k2^{k/2} \log\left(\frac{n}{k}\right)\right)$
Approx min-wise, uniform distrib on $\mathcal{F}$	$O(n^2/\epsilon^2)$ (existential) ??? (constructive)	$n^2(1 - \sqrt{8\epsilon})$
Approx min-wise, biased distrib on $\mathcal{F}$	???	$\max_{r \geq 1} \frac{(n-r) \binom{n}{r}}{1 + \epsilon \binom{n}{r}}$
Approx min-wise, restricted, uniform distrib on $\mathcal{F}$	$O\left(\frac{k^2 \log(n/k)}{\epsilon^2}\right)$ (existential) $2^{4k+o(k)} k^{2 \log(\log n/\epsilon)}$ (constructive)	?
Approx min-wise, restricted, biased distrib on $\mathcal{F}$	?	$\Omega\left(\min\left(\frac{\log(n/k)}{\epsilon^{1/3}}, k2^{k/2} \log(n/k)\right)\right)$

Table 1: Synopsis of results – minimum size of families

Family type	Bounds on the most probable element		Bounds on the least probable element	
	Upper	Lower	Upper	Lower
Pairwise independent – worst set	$O\left(\frac{1}{\sqrt{k}}\right)$	?	???	$\frac{1}{2(k-1)}$
Linear – worst set	?	$\frac{3 \ln k}{\pi^2 k}$	$\frac{12 \ln 2}{\pi^2 k}$	?
Pairwise independent – random set	$\frac{1 + 1/\sqrt{2}}{k}$	???	???	?
Linear – random set	?	???	???	?

Table 2: Synopsis of results – quality of approximation

## 2.1 Exact problem with non-uniform distribution

Although we focus on results for uniform distributions, we demonstrate here an interesting result: the lower bound of Theorem 1 can be beaten by using non-uniform distributions.

**Theorem 3** *There is a family  $\mathcal{F}$  of size at most  $n2^{n-1} - 1$ , such that  $\mathcal{F}$  with an associated distribution  $\mu$  is min-wise independent.*

*Proof:* We can write a linear program to find a  $\mathcal{F}$  and  $\mu$  satisfying the theorem. We have a variable  $x_\pi$  for each of the

permutations  $\pi \in S_n$ , where  $x_\pi$  represents the weight of  $\pi$  according to  $\mu$ . For every  $X \subset [n]$  and for every  $x \in X$ , we express the condition that  $\Pr(\min\{\pi(X)\} = \pi(x)) = \frac{1}{|X|}$  as a linear equation in the variables  $x_\pi$ . We have a total of  $\sum_{k=1}^n k \cdot \binom{n}{k} = n2^{n-1} - 1$  constraints. This system clearly has a feasible solution (choose an element of  $S_n$  uniformly at random, that is, put  $x_\pi = 1/n!$  for all  $\pi \in S_n$ ), and hence it has a basic feasible solution with at most  $n \cdot 2^{n-1} - 1$  non-zero variables. This solution yields a family satisfying the conditions of the theorem.  $\square$

**Remark 2** Although this beats the lower bound of Corollary 1, the size of the family is still exponential in  $n$ , and we will prove an almost tight lower bound in Section 3.3. Also, for restricted min-wise independence, this same construction gives an upper bound of  $\sum_{j=1}^k j \cdot \binom{n}{j}$ .

### 3 The Approximate Problem

As the exact problem requires exponential sized families, we turn our attention to the approximate problem.

#### 3.1 Existential Upper Bounds

We can obtain existential upper bounds on the sizes of approximately min-wise independent families via the probabilistic method [2], by simply choosing a number of random permutations from  $S_n$ . We omit the straightforward proofs.

**Theorem 4** *There exist families of size  $O(\frac{n^2}{\epsilon^2})$  that are approximately min-wise independent and there exist families of size  $O(\frac{k^2 \log(n/k)}{\epsilon^2})$  that are approximately and restricted min-wise independent.*

In fact, a family of permutations of size  $O(\frac{n^2}{\epsilon^2})$  chosen uniformly at random from  $S_n$  will be approximately min-wise independent with high probability. This would appear to provide a suitable solution for the document similarity problem discussed in the introduction. In practice, however, this result does not help us, since one cannot conveniently represent a random permutation from  $S_n$ . Recall that a random permutation on  $n$  elements requires on average  $\Omega(n \log n)$  bits to represent, and in practice  $n = 2^{64}$ . This leads us to consider simple linear permutations in Section 4.

#### 3.2 Lower Bound for Uniform Families

We will prove a lower bound of  $n^2(1 - \sqrt{8\epsilon})$  for families with the uniform probability distribution. This shows that the  $n^2$  term in the existential upper bound cannot be improved.

**Theorem 5** *Let  $\mathcal{F}$  be an approximate min-wise independent family. Then  $|\mathcal{F}| \geq n^2(1 - \sqrt{8\epsilon})$ .*

*Proof:* Let  $|\mathcal{F}| = f$ . There must be some element  $a$  such that  $\pi(a) = 1$  (that is,  $a$  is the second smallest after the permutation) for at least  $f/n$  permutations of  $\mathcal{F}$ . Fix such an  $a$  and consider  $z \leq f/n$  such permutations. We will choose a value for  $z$  later. Let  $Z$  be the set of elements which occur as the smallest element in these  $z$  permutations (that is,  $b \in Z$  iff  $\pi(b) = 0$  for at least one of these  $z$  permutations) and let  $S = [n] - Z$ . Clearly  $a \in S$  and  $|S| \geq n - z$ . Consider for how many permutations  $\pi \in \mathcal{F}$  it is the case that  $\pi(a)$  is the smallest element of  $\pi(S)$ . This happens at least whenever  $\pi(a) = 0$  and also for the  $z$  permutations discussed above, where  $\pi(a) = 1$  but an element not in  $S$  has image 0 under  $\pi$ . But  $\pi(a) = 0$  for at least  $\frac{f}{n}(1 - \epsilon)$  permutations, because  $\mathcal{F}$  is an approximately min-wise independent family; and for

the same reason,  $\pi(a)$  can be the minimum element of  $S$  for at most  $\frac{f}{|S|}(1 + \epsilon) \leq \frac{f(1+\epsilon)}{n-z}$  permutations. Hence

$$\frac{f(1 - \epsilon)}{n} + z \leq \frac{f(1 + \epsilon)}{n - z}.$$

Solving this equation for  $f$  and (almost) optimizing for  $z$  ( $z = \sqrt{2\epsilon}f/n$ ) yields

$$f \geq n^2 \frac{1 - \sqrt{2\epsilon}}{1 + \sqrt{2\epsilon} - \epsilon}.$$

Simplifying the above yields a lower bound of  $n^2(1 - \sqrt{8\epsilon})$  on  $|\mathcal{F}|$ .  $\square$

#### 3.3 Lower Bound for Non-Uniform Families

We will prove a lower bound on the size of any approximately min-wise independent family, even non-uniform families with an associated probability distribution  $\mu$ . Our lower bound proof also yields a lower bound for non-uniform exactly min-wise independent families that is very close to the upper bound of  $n2^{n-1} - 1$  obtained in Section 2.1.

**Theorem 6** *Let  $\mathcal{F}$  be an approximate min-wise independent family, possibly with an associated probability distribution  $\mu$ . Then  $|\mathcal{F}| \geq \frac{(n-r)\binom{n}{r}}{1 + \epsilon 2^r \binom{n}{r}}$ , for any  $r < n$ .*

*Proof:* Fix an element  $a$  and a set  $Z = \{x_1, x_2, \dots, x_r\} \subseteq [n]$  with  $a \notin Z$ . Let us say that the pair  $(Z, a)$  is *satisfied* if there is a permutation  $\pi$  in  $\mathcal{F}$  that has all the elements of  $\pi(Z)$  as the  $r$  smallest elements of  $\pi$  in any order (that is,  $\pi(Z) = [r]$ ) and has  $a$  as the  $(r+1)$ st smallest element (that is,  $\pi(a) = r+1$ ). We will show that most pairs  $(Z, a)$  must be satisfied for  $\mathcal{F}$  to be an approximately min-wise independent family, and that in fact all pairs  $(Z, a)$  must be satisfied for  $\mathcal{F}$  to be an exactly min-wise independent family.

Let  $Y = [n] - Z$ . By definition  $a \in Y$ . We consider the sets  $Y_i = Y \cup x_i$  and count how often  $\pi(a)$  is the smallest element of  $\pi(Y_i)$ . Let  $B_S$  be the event that  $a$  is the minimum of  $\pi(S)$  when we choose a permutation from  $\mathcal{F}$  under the distribution  $\mu$ . Let  $B = \bigcup_{i=1}^r B_{Y_i}$ . Then  $B \subseteq B_Y$ , and hence  $\Pr(B_Y - B) = \Pr(B_Y) - \Pr(B)$ . On the other hand, the event  $B_Y - B$  is precisely the event that  $(Z, a)$  is satisfied.

We now use the inclusion-exclusion principle to calculate  $\Pr(B) = \Pr(\bigcup_{i=1}^r B_{Y_i})$ . It is helpful to note the following facts. First if  $a \in S_2 \subseteq S_1$  then  $B_{S_1} \subseteq B_{S_2}$  and if  $a \in S_1 \cap S_2$  then  $B_{S_1} \cap B_{S_2} = B_{S_1 \cup S_2}$ . Second, by the definition of approximate min-wise independence,  $\frac{1-\epsilon}{|S|} \leq \Pr(B_S) \leq \frac{1+\epsilon}{|S|}$ . We will abbreviate this by saying that  $\Pr(B_S) = \frac{1 \pm \epsilon}{|S|}$ , where the meaning is clear. Third, the union of  $i$  distinct  $Y_i$ 's has size  $n - r + i$ . Hence

$$\begin{aligned} \Pr(B) &= \Pr(B_{Y_1}) + \Pr(B_{Y_2}) + \dots + \Pr(B_{Y_n}) \\ &\quad - \Pr(B_{Y_1} \cap B_{Y_2}) - \dots \end{aligned}$$

$$\begin{aligned}
& + \Pr(B_{Y_1} \cap B_{Y_2} \cap B_{Y_3}) + \dots \\
& = \Pr(B_{Y_1}) + \Pr(B_{Y_2}) + \dots + \Pr(B_{Y_n}) \\
& \quad - \Pr(B_{Y_1 \cup Y_2}) - \dots + \Pr(B_{Y_1 \cup Y_2 \cup Y_3}) + \dots \\
& = \sum_{i=1}^r (-1)^{i+1} \binom{r}{i} \frac{1 \pm \epsilon}{n-r+i}
\end{aligned}$$

Hence

$$\begin{aligned}
& \Pr(B_Y - B) \\
& = \frac{1 \pm \epsilon}{n-r} - \sum_{i=1}^r (-1)^{i+1} \binom{r}{i} \frac{1 \pm \epsilon}{n-r+i} \\
& = \sum_{i=0}^r (-1)^i \binom{r}{i} \frac{1 \pm \epsilon}{n-r+i} \\
& = \sum_{i=0}^r (-1)^i \binom{r}{i} \frac{1}{n-r+i} \pm \epsilon \sum_{i=0}^r \binom{r}{i} \frac{1}{n-r+i}
\end{aligned}$$

To evaluate the first term in the expression above, note that it equals  $\Pr(B_Y - B)$  when  $\epsilon$  is 0. That is, the term is the probability that  $(Z, a)$  is satisfied for an exactly min-wise independent family. Note that it depends only on  $n$  and  $r$ , and not on the family under consideration! In particular, we can calculate it easily by computing the probability that  $(Z, a)$  is satisfied for the family  $S_n$ , which is  $\frac{1}{\binom{n-r}{r}}$ . (Thus we obtain the combinatorial identity

$$\sum_{i=0}^r (-1)^i \binom{r}{i} \frac{1}{n-r+i} = \frac{1}{\binom{n-r}{r}}.$$

The hint for its algebraic derivation is [10, equation 1.2.6.24].)

The magnitude of the coefficient of  $\epsilon$  is at most  $\frac{2^r}{n-r}$ . Hence

$$\begin{aligned}
& \frac{1}{\binom{n-r}{r}} + \epsilon \frac{2^r}{n-r} \\
& \geq \Pr(B_Y - B) \geq \\
& \quad \frac{1}{\binom{n-r}{r}} - \epsilon \frac{2^r}{n-r}
\end{aligned} \tag{8}$$

Since  $\Pr(B_Y - B) \leq \frac{1}{\binom{n-r}{r}} + \epsilon \frac{2^r}{n-r}$ , the total probability mass of the permutations that satisfy any given pair  $(Z, a)$  is at most  $p = \frac{1}{\binom{n-r}{r}} + \epsilon \frac{2^r}{n-r}$ . Hence the number of distinct pairs  $(Z, a)$  which have some permutation satisfying them must be at least  $1/p$ . But every permutation satisfies exactly one  $(Z, a)$  pair. This means that there must be at least  $1/p$  permutations, that is, the size of the family is at least  $\frac{\binom{n-r}{r}}{1 + \epsilon 2^r}$ .  $\square$

**Corollary 1** *Let  $\mathcal{F}$  be exact min-wise independent family, possibly with an associated probability distribution  $\mu$ . Then  $|\mathcal{F}| \geq \lceil \frac{n}{2} \rceil \binom{n}{\lfloor n/2 \rfloor}$ .*

*Proof:* Plug  $\epsilon = 0$  and  $r = \lfloor \frac{n}{2} \rfloor$  in the result of Theorem 6  $\square$

Actually, Theorem 6 proves an even stronger corollary: Equation (8) shows that the probability that  $(Z, a)$  is satisfied is positive as long as  $\epsilon < 1/2^r \binom{n}{r}$ . Hence, for any approximate min-wise independent family with such an  $\epsilon$ , all  $\binom{n}{r} (n-r)$  possible pairs  $(Z, a)$  are satisfied, and hence there are at least this many permutations. This is maximized for  $r = \lfloor \frac{n}{2} \rfloor$ , and hence the bound of Corollary 1 also holds for approximate families with an exponentially small  $\epsilon$ .

## 4 Linear and Pairwise Independent Families

We now focus on the behavior of permutations most likely to be used in practice, linear transformations. In particular, we focus on the situation where the universe of elements is  $[p]$  for some prime  $p$ , and the family of permutations is given by all permutations of the form  $\pi(x) = ax + b \pmod{p}$  (with  $a \neq 0$ ). Linear transformations are easy to represent and efficiently calculable, making them suitable for real applications. Our results suggest that although this family of permutations is not min-wise independent, its performance should be sufficient in many practical situations.

### 4.1 General Upper and Lower Bounds

As the results for linear permutations require significant calculations, we do not provide proofs for all the results here. We begin with a simple lower bound that holds not just for linear transformations but for any pairwise independent family of permutations; many of our results have this form.

**Theorem 7** *For any  $X \subseteq [n]$  with  $|X| = k$  and for any  $x \in X$ ,*

$$\Pr(\min\{\pi(X)\} = \pi(x)) > \frac{1}{2(k-1)}$$

*if  $\pi$  is chosen from a pairwise independent family of permutations.*

*Proof:* Consider a set  $X = \{x_0, \dots, x_{k-1}\}$ . We will show that  $\pi(x_0)$  is the smallest element of  $\pi(X)$  as often as required by the theorem. Suppose that  $\pi(x_0) = z$ . If  $\pi$  is chosen from a pairwise independent family, then  $\Pr(\pi(x_i) < z | \pi(x_0) = z) = z/n$ . Since the probability that  $\pi$  maps  $x_i$  to something smaller than  $\pi(x_0)$  is  $z/n$ , the probability that  $\pi$  maps any element of  $X$  to something smaller than  $\pi(x_0)$  is at most  $(k-1)z/n$ , and hence  $\pi(x_0)$  is the minimum of  $\pi(X)$  with probability at least  $1 - (k-1)z/n$ . This is non-negative for  $0 \leq z \leq \lfloor \frac{n}{k-1} \rfloor$ . Hence

$$\begin{aligned}
\Pr(\min\{\pi(X)\} = \pi(x_0)) & \geq \frac{1}{n} \sum_{z=0}^{\lfloor n/(k-1) \rfloor} \left(1 - \frac{(k-1)z}{n}\right) \\
& > \frac{1}{2(k-1)}
\end{aligned}$$

$\square$

We have an upper bound on  $\Pr(\min\{\pi(X)\} = \pi(x))$  for all pairwise independent families of permutations that is  $O(1/\sqrt{k})$ , based on a linear programming formulation of the problem. Subsequent to our original proof, Piotr Indyk suggested a simpler proof for this bound.

## 4.2 Linear Families, Upper and Lower Bounds

We can derive further bounds by considering specifically linear transformations. For instance, we can show that the family of linear transformations is not even approximately minwise independent for any constant  $\epsilon$ . Here we sketch this result.

**Theorem 8** *Consider the set  $X_k = \{0, 1, 2 \dots k\}$ , as a subset of  $[p]$ . As  $k, p \rightarrow \infty$ , with  $p \gg k$ ,*

$$\Pr(\min\{\pi(X)\} = \pi(0)) \sim \frac{3 \ln k}{\pi^2 k}$$

when  $\pi$  is a randomly chosen linear transformation of the form  $\pi(x) = ax + b \bmod p$  (with  $a \neq 0$ ).

*Proof:* The proof will use some basic facts about *Farey series*. We first remind the reader of the definition and some basic facts regarding Farey series; more information can be found in most standard number theory texts.

**Definition 1** *The Farey series of order  $k$  consists of all irreducible fractions less than 1 with denominator at most  $k$ , in increasing order.*

If  $\frac{n_1}{d_1}$  and  $\frac{n_2}{d_2}$  are two consecutive fractions in the order  $k$  Farey series then

1.  $n_2 d_1 - n_1 d_2 = 1$ .
2.  $(d_1, d_2) = 1$ .
3. the first fraction inserted between  $\frac{n_1}{d_1}$  and  $\frac{n_2}{d_2}$  in a higher order Farey series is  $\frac{n_1+n_2}{d_1+d_2}$ .

To compute the fraction of time that  $\pi(0)$  is the minimum element of  $\{\pi(X_k)\}$ , let us first consider all transformations  $\pi$  with multiplier  $a$ . Let  $z_a = \min_{i=1, \dots, k} \{-a \cdot i \bmod p\}$ . Then  $\pi(0)$  is minimal only for those  $\pi = ax + b \bmod p$  where  $b < z_a$  (note that  $z_a$  is positive!), since for the other values of  $b$  the image of the minimal element will lie behind  $\pi(0) = b$ .

Hence, to find the fraction of the time that 0 is the minimum element of  $\{\pi(X_k)\}$ , it suffices to find the expected value of  $\frac{1}{p} \min_{i=1, \dots, k} \{-a \cdot i \bmod p\}$ , which conveniently is also the expected value of  $\frac{1}{p} \min_{i=1, \dots, k} \{a \cdot i \bmod p\}$ . We concentrate on the latter expression.

Consider what happens to the numbers  $\{a \cdot i \bmod p | i = 1 \dots k\}$  as we increase the value of the multiplier  $a$  from 1 to  $p-1$ . It is useful to think of the numbers  $0, \dots, p-1$  as arranged clockwise around a circle. Consider  $k$  tokens, corresponding to the numbers  $1, \dots, k$  from the set  $X_k$ . For

each  $i$ , we view  $a \cdot i \bmod p$  as the position of the  $i$ th token at time  $a$ . Token  $i$  starts in position  $i$ . As we increase the value of the multiplier  $a$  from 1 to  $p-1$  all tokens move around the circle in clockwise direction but at different speeds: token  $i$  moves  $i$  steps for every time tick.

If  $p$  is sufficiently larger than  $k$ , we can think of this motion as being continuous. That is, we scale the circle so that its circumference is 1. Let  $f = \frac{a}{p}$ . Then the distance of token  $i$  from the origin along the circle when the multiplier is  $a$  is the fractional part of  $f \cdot i$ . Henceforth we think of this motion of the tokens as being continuous, with the “time”  $f$  increasing uniformly from 0 to 1. We need to compute the average distance of the token closest to the origin as  $f$  increases uniformly from 0 to 1, where distance here is measured as clockwise distance along the circumference. This average distance is (asymptotically)  $\frac{1}{p} \min_{i=1, \dots, k} \{a \cdot i \bmod p\}$ , the term we wish to compute. (Asymptotically this approximation yields the correct answer, as the approximation affects only lower order terms.)

The token closest to the origin changes whenever a token reaches the origin. This happens whenever the value of  $f$  is  $\frac{n}{d}$  for integers  $n$  and  $d$  with  $1 \leq n < d \leq k$ , as at that point the token with speed  $d$  reaches the origin. Thus the times where the token closest to the origin changes are precisely the proper (less than 1) fractions of denominator at most  $k$ , that is, the terms of the Farey sequence of order  $k$ . Let  $\frac{n_1}{d_1}$  and  $\frac{n_2}{d_2}$  be two consecutive fractions in the Farey sequence of order  $k$ . For  $\frac{n_1}{d_1} \leq f \leq \frac{n_2}{d_2}$ , the token with speed  $d_1$  is closest to the origin. This time interval has length  $\frac{n_2}{d_2} - \frac{n_1}{d_1} = \frac{1}{d_1 d_2}$ . During this time interval, the token starts at the origin and moves with a speed of  $d_1$ . Thus the average distance of this token from the origin during this interval is  $\frac{1}{2} \cdot d_1 \cdot \frac{1}{d_1 d_2} = \frac{1}{2 d_2}$ .

To obtain the average distance over the entire interval, it suffices to take the appropriate weighted sum over all pairs of consecutive Farey fractions. By the above, the contribution from each interval  $[\frac{n_1}{d_1}, \frac{n_2}{d_2}]$  is  $\frac{1}{d_1 d_2} \cdot \frac{1}{2 d_2} = \frac{1}{2 d_1 d_2^2}$ .

To find a simple form for the resulting sum, we build up, starting the appropriate sum for  $X_1 = \{0, 1\}$  and building up to the set  $X_k$ . Alternatively, we may think of how the sum changes as we build up from the order  $j-1$  Farey series to the order  $j$  Farey series and use this to derive the appropriate sum for the order  $k$  Farey series. The order  $j$  Farey series is derived from the order  $j-1$  Farey series by adding all fractions of the form  $\frac{a}{j}$  with  $(a, j) = 1$  in their proper position. (Note we use the standard shorthand  $(a, j)$  for  $\gcd(a, j)$ .) Correspondingly, this changes the contribution to the summation in all intervals where a new fraction is inserted. Suppose a fraction is inserted between  $\frac{n_1}{d_1}$  and  $\frac{n_2}{d_2}$ . Then the inserted fraction must be  $\frac{n_1+n_2}{d_1+d_2}$ , where  $d_1+d_2 = k$ . Before the insertion, the contribution of this interval was  $\frac{1}{2 d_1 d_2^2}$ . After the insertion, the contribution becomes  $\frac{1}{2 d_1 (d_1+d_2)^2} + \frac{1}{2 (d_1+d_2) d_2^2}$ . Thus the change is

$$\frac{1}{2 d_1 (d_1 + d_2)^2} + \frac{1}{2 (d_1 + d_2) d_2^2} - \frac{1}{2 d_1 d_2^2}$$

$$\begin{aligned}
&= \frac{d_2^2 + d_1(d_1 + d_2) - (d_1 + d_2)^2}{2d_1(d_1 + d_2)^2 d_2^2} \\
&= \frac{-1}{2(d_1 + d_2)^2 d_2}
\end{aligned}$$

Note that  $d_1 + d_2 = j$ . Further  $(j, d_2) = 1$ . In fact, for every  $a$  such that  $(a, j) = 1$ , there exists two consecutive Farey fractions  $\frac{a_1}{d_1}$  and  $\frac{a_2}{d_2}$  such that  $d_1 + d_2 = j$  and  $d_2 = a$ . Thus the change in the summation caused by building up from order  $j-1$  to order  $j$  Farey sequences is  $\frac{1}{2j^2} \sum_{(a,j)=1, 1 \leq a \leq j} \frac{1}{a}$ . For the order 1 Farey sequence, the value of the appropriate summation is obviously  $\frac{1}{2}$ . Thus the value for the order  $k$  Farey sequence is

$$\frac{1}{2} \left( 1 - \sum_{j=2}^k \frac{1}{j^2} \sum_{(a,j)=1, 1 \leq a \leq j} \frac{1}{a} \right)$$

From here one must simply evaluate the value of this expression asymptotically to obtain the theorem. We spare the reader the apparently unenlightening algebraic details.  $\square$

Theorem 8 shows that it is possible to find sets for which some element is minimal for  $\Omega(\frac{\log k}{k})$  of the time when random linear transformations are used. Similarly, we can show that  $\pi(\lfloor (k-1)/2 \rfloor)$  is asymptotically the minimum element of  $\pi(X_k)$  with probability  $\frac{12 \ln 2}{\pi^2 k} \approx \frac{0.843}{k}$ . This result provides an example of how much less often than  $\frac{1}{k}$  of the time an element can be minimal when random linear transformations are used.

Despite the seemingly bad worst-case behavior of linear transformations, we believe that in practice they are suitable for applications, because they perform well on random sets. For a set  $X = \{x_0, \dots, x_{k-1}\}$  of size  $k$ , let  $F(X)$  be  $\max_i \frac{|\{\pi \mid \min\{\pi(X)\} = \pi(x_i)\}|}{p(p-1)}$ . That is,  $F(X)$  is the fraction of the permutations for which the most likely element to be the minimum is actually the minimum. (And we have just seen that  $F(X) \geq \frac{3}{\pi^2} \frac{\ln k}{k}$  in the worst case.) We now prove that the expected value of  $F(X)$  when  $X$  is chosen uniformly at random from all sets of size  $k$  as  $k, p \rightarrow \infty$  can be bounded by  $(1 + 1/\sqrt{2})/k + O(1/k^2)$ . In this sense, linear transformations are approximately min-wise independent with respect to random sets.

**Theorem 9** *As  $k, p \rightarrow \infty$ , with  $p \gg k^2$ ,  $\mathbf{E}_X[F(X)]$  is bounded above by  $(1 + 1/\sqrt{2})/k + O(1/k^2)$ .*

*Proof:* We define

$$f_i(X) = \frac{|\{\pi \mid \min\{\pi(X)\} = \pi(x_i)\}|}{p(p-1)},$$

and

$$g_i(z, X) = \frac{|\{\pi \mid \min\{\pi(X)\} = \pi(x_i) \text{ and } \pi(x_i) = z \cdot p\}|}{p-1},$$

That is, consider the subset of permutations that map the  $i$ th element to  $zp$ . Then  $g_i$  is the fraction of these permutations for which the  $i$ th element is minimal.

Hereafter we suppose that the universe size  $p$  is sufficiently large that we may think of  $z$  as varying continuously on the unit circle from 0 to 1, instead of jumping discretely by  $1/p$ . This simplification allows us to dismiss many lower order terms. Similarly, we will suppose that  $p$  is sufficiently large compared to  $k$  so that we may suppose that the  $k$  values of  $X$  are chosen with replacement, and the results will be equivalent asymptotically.

The value we wish to bound is

$$F(X) = \mathbf{E}_X \left[ \max_{i=0, \dots, k-1} f_i(X) \right],$$

where we use  $\mathbf{E}_X$  to denote that the expectation is over the random choice of the set  $X$ . Note also that we have the following relation:

$$f_i(X) = \int_0^1 g_i(z, X) dz.$$

Let the  $f_i(X)$  have mean  $\mu$  and variance  $\sigma^2$ . (Note the mean and variance are the same for all  $f_i$ .) To bound  $F(X)$ , we make use of a simple bound on the expected value of the maximum of several identically distributed random variables.

**Lemma 1** *Let  $X_1, X_2, \dots, X_k$  be identically distributed random variables with mean  $\mu$  and variance  $\sigma^2$ . Then*

$$\mathbf{E} \left[ \max_{i=1, \dots, k} X_i \right] \leq \mu + \sigma \sqrt{k}.$$

*Proof:* We show equivalently that

$$\left( \mathbf{E} \left[ \max_{i=1, \dots, k} X_i - \mu \right] \right)^2 \leq k \sigma^2.$$

$$\begin{aligned}
\left( \mathbf{E} \left[ \max_{i=1, \dots, k} X_i - \mu \right] \right)^2 &\leq \mathbf{E} \left[ \left( \max_{i=1, \dots, k} X_i - \mu \right)^2 \right] \\
&\leq \mathbf{E} \left[ \max_{i=1, \dots, k} (X_i - \mu)^2 \right] \\
&\leq \mathbf{E} \left[ \sum_{i=1, \dots, k} (X_i - \mu)^2 \right] \\
&\leq \sum_{i=1, \dots, k} \mathbf{E} [(X_i - \mu)^2] \\
&= k \sigma^2
\end{aligned}$$

$\square$

Clearly, by symmetry  $\mathbf{E}_X[f_i(X)] = 1/k$ . Hence, to find an upper bound on  $F$ , we just have to bound  $\sigma^2$ , the variance of  $f_i(X)$ . Specifically, we bound the variance of  $f_0(X)$ .

We define some helpful notation. Let  $\pi_{a,z}$  denote the unique linear permutation such that  $ax_0 + b = z \cdot p \pmod p$ . That is,  $\pi_{a,z}$  is the linear permutation with multiplier  $a$

that maps  $x_0$  to  $z \cdot p$ . Let  $M_a(z, X)$  be an indicator random variable that is 1 if  $\min\{\pi_{a,z}(X)\} = \pi_{a,z}(x_0)$ . Thus,  $g_0(z, X) = \frac{1}{p} \sum_a M_a(z, X)$ . Now the variance of  $f_0$  is just

$$\begin{aligned} \sigma^2 &= \mathbf{E}_X [(f_0(X) - \mathbf{E}_X[f_0(X)])^2] \\ &= \mathbf{E}_X \left[ \left( \int_0^1 g_0(z, X) dz - \mathbf{E}_X \left[ \int_0^1 g_0(z, X) dz \right] \right)^2 \right] \\ &= \mathbf{E}_X \left[ \left( \int_0^1 (g_0(z, X) - \mathbf{E}_X[g_0(z, X)]) dz \right)^2 \right] \\ &= \frac{1}{p^2} \mathbf{E}_X \left[ \left( \int_0^1 \left( \sum_a M_a(z, X) - \mathbf{E}_X \left[ \sum_a M_a(z, X) \right] \right) dz \right)^2 \right] \\ &= \frac{1}{p^2} \mathbf{E}_X \left[ \left( \int_0^1 \sum_a \left( M_a(z, X) - \mathbf{E}_X[M_a(z, X)] \right) dz \right)^2 \right] \end{aligned}$$

Let  $\mu_a(z) = \mathbf{E}_X[M_a(z, X)]$ . From this definition, it is apparent that  $\mu_a(z) = (1-z)^{k-1}$ , as each of the images of the other randomly chosen  $k-1$  elements has probability  $1-z$  of being greater than  $z \cdot p$ .

Hence, continuing from the last line above,

$$\begin{aligned} \sigma^2 &= \frac{1}{p^2} \mathbf{E}_X \left[ \left( \int_0^1 \sum_a (M_a(z, X) - \mathbf{E}_X[M_a(z, X)]) dz \right)^2 \right] \\ &= \frac{1}{p^2} \mathbf{E}_X \left[ \int_{z=0}^1 \int_{y=0}^1 \left( \sum_{a_1, a_2} (M_{a_1}(z, X) - \mu_{a_1}(z)) \right. \right. \\ &\quad \left. \left. \times (M_{a_2}(y, X) - \mu_{a_2}(y)) \right) dy dz \right] \\ &= \frac{1}{p^2} \int_{z=0}^1 \int_{y=0}^1 \left( \sum_{a_1, a_2} \left( \mathbf{E}_X[M_{a_1}(z, X)M_{a_2}(y, X)] \right. \right. \\ &\quad \left. \left. - \mu_{a_1}(z)\mu_{a_2}(y) \right) \right) dy dz \end{aligned} \tag{9}$$

We now bound the last term. This will in turn bound the variance and yield the theorem. In order to do this, we derive an alternative expression for  $\mathbf{E}_X[M_{a_1}(z, X)M_{a_2}(y, X)]$  that can be appropriately bounded.

Let

$$q_{a_1, a_2}(z, y) = \Pr_{x \in [p]} (\pi_{a_1, z}(x) > y \cdot p \text{ and } \pi_{a_2, y}(x) > z \cdot p).$$

Then

$$\mathbf{E}_X[M_{a_1}(z, X)M_{a_2}(y, X)] = (q_{a_1, a_2}(z, y))^{k-1},$$

again since the other  $k-1$  terms of  $X$  are chosen uniformly at random.

We thus have expressed the value we wish to bound as the sum of the  $(k-1)$ st powers of  $q_{a_1, a_2}$  terms. The next lemma shows that the sum of these  $q_{a_1, a_2}$  terms is fixed. As the maximum possible value of the sum of the  $(k-1)$ st powers is achieved when the terms in the sum take on extremal values, together these results will allow us to bound  $\sum_{a_1, a_2} \mathbf{E}_X[M_{a_1}(z, X)M_{a_2}(y, X)]$ .

**Lemma 2**

$$\sum_{a_1, a_2} q_{a_1, a_2}(z, y) = p^2(1-z)(1-y).$$

*Proof:* Consider the following experiment. We choose three values  $a_1, a_2, x \in [p]$  independently and uniformly at random. The experiment succeeds if both  $\pi_{a_1, z}(x) > y \cdot p$  and  $\pi_{a_2, y}(x) > z \cdot p$ . Clearly, the probability of success is  $(1-z)(1-y)$ . The summation  $\sum_{a_1, a_2} p \cdot q_{a_1, a_2}(z, y)$  is simply the number of the  $p^3$  triples  $(a_1, a_2, x)$  for which the experiment succeeds. The lemma follows.  $\square$

Since the total sum of the terms  $q_{a_1, a_2}$  is fixed, the sum  $\sum_{a_1, a_2} \mathbf{E}_X[M_{a_1}(z, X)M_{a_2}(y, X)]$  is maximized when the  $q_{a_1, a_2}$  terms take on extremal values. Let us assume that  $z \geq y$ . Then  $q_{a_1, a_2}(z, y) \in [1-z-y, 1-z]$ . (Of course  $q_{a_1, a_2}(z, y) \geq 0$ , and hence the above range may not be correct if  $z+y > 1$ .) A simple calculation then yields the following bound (for  $z+y \leq 1$ ):

$$\begin{aligned} \sum_{a_1, a_2} \mathbf{E}_X[M_{a_1}(z, X)M_{a_2}(y, X)] \\ \leq p^2 (z(1-z)^{k-1} + (1-z)(1-z-y)^{k-1}). \end{aligned}$$

We will use this bound for the range  $z \leq 1/2$ . For  $z > 1/2$ , we have  $q_{a_1, a_2}(z, y) \leq 1-z \leq 1/2$ . Hence,

$$\sum_{a_1, a_2} \mathbf{E}_X[M_{a_1}(z, X)M_{a_2}(y, X)] \leq p^2(1/2^{k-1}).$$

Substituting this bound in (9), we get:

$$\begin{aligned} \sigma^2 &= \frac{1}{p^2} \mathbf{E}_X \left[ \int_{z=0}^1 \int_{y=0}^1 \left( \sum_{a_1, a_2} (M_{a_1}(z, X)M_{a_2}(y, X) \right. \right. \\ &\quad \left. \left. - \mu_{a_1}(z)\mu_{a_2}(y)) \right) dy dz \right] \\ &= \frac{2}{p^2} \int_{z=0}^1 \int_{y=0}^z \left( \sum_{a_1, a_2} \mathbf{E}_X \left[ M_{a_1}(z, X)M_{a_2}(y, X) \right. \right. \\ &\quad \left. \left. - \mu_{a_1}(z)\mu_{a_2}(y) \right] \right) dy dz \\ &\leq 2 \int_{z=0}^{1/2} \int_{y=0}^z \left( z(1-z)^{k-1} + (1-z)(1-z-y)^{k-1} \right. \\ &\quad \left. - (1-z)^{k-1}(1-y)^{k-1} \right) dy dz \\ &\quad + 2 \int_{z=1/2}^1 \int_{y=0}^z \frac{1}{2^{k-1}} dy dz \end{aligned}$$

To prove Theorem 1, we need merely to compute this integral thus bounding the variance. This calculation is easily performed, yielding

$$\sigma^2 \leq \frac{1}{2k^3} + O(1/k^4).$$

This proves Theorem 9.  $\square$

Simulations suggest that in fact the behavior of families of linear transformations on a random set  $X$  is much better than this. We conjecture that the expected value of  $F(X)$  converges to  $1/k$  asymptotically.

Also, we note that Theorem 9 actually generalizes quite straightforwardly to all pairwise independent families. The notation becomes slightly more difficult, but the proof follows the same course.

## 5 Acknowledgments

The authors thank Noam Elkies for enlightening discussions regarding Farey series.

## References

- [1] N. Alon, M. Dietzfelbinger, P. B. Miltersen, E. Petrank, and G. Tardos. Is linear hashing good? In *Proceedings of the Twenty-Ninth Annual ACM Symposium on Theory of Computing*, pages 465–474, El Paso, Texas, 4–6 May 1997.
- [2] N. Alon and J. H. Spencer. *The Probabilistic Method*. John Wiley and Sons, 1992.
- [3] T. M. Apostol. *Introduction to Analytic Number Theory*. Springer-Verlag, 1976.
- [4] T. Berners-Lee, R. Cailliau, A. Loutonen, H. F. Nielsen, and A. Secret. The world-wide web. *Communications of the ACM*, 37(8):76–82, 1994.
- [5] A. Z. Broder. On the resemblance and containment of documents. In *Proceedings of Compression and Complexity of SEQUENCES 1997*. To appear.
- [6] A. Z. Broder, S. C. Glassman, M. S. Manasse, and G. Zweig. Syntactic clustering of the Web. In *Proceedings of the Sixth International World Wide Web Conference*, pages 391–404, 1997.
- [7] A. Z. Broder and A. R. Karlin. Multilevel adaptive hashing. In *Proceedings of the First Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 43–53, San Francisco, California, 22–24 Jan. 1990.
- [8] J. L. Carter and M. N. Wegman. Universal classes of hash functions. *Journal of Computer and System Sciences*, 18(2):143–154, Apr. 1979.
- [9] M. Dietzfelbinger, A. Karlin, K. Mehlhorn, F. M. auf der Heide, H. Rohnert, and R. E. Tarjan. Dynamic perfect hashing: Upper and lower bounds. *SIAM J. Comput.*, 23(4):738–761, Aug. 1994.
- [10] D. E. Knuth. *The Art of Computer Programming, Vol. I: Fundamental Algorithms*. Addison-Wesley, second edition, 1973.
- [11] M. Luby and A. Wigderson. Pairwise independence and derandomization. Technical Report TR-95-035, International Computer Science Institute, Berkeley, California, 1995.
- [12] R. Seltzer, E. J. Ray, and D. S. Ray. *The Alta Vista Search Revolution : How to Find Anything on the Internet*. McGraw-Hill, 1996.
- [13] R. J. Souza, P. Krishnakumar, C. M. Özveren, R. J. Simcoe, B. A. Spinney, R. E. Thomas, and R. J. Walsh. GIGAswitch: A high-performance packet-switching platform. *DIGITAL Technical Journal*, 6(1):9–22, 1994.