

Linear Hashing is Awesome

Mathias Bæk Tejs Knudsen*

University of Copenhagen,
mathias@tejs.dk

Abstract

We consider the hash function $h(x) = ((ax + b) \bmod p) \bmod n$ where a, b are chosen uniformly at random from $\{0, 1, \dots, p - 1\}$. We prove that when we use $h(x)$ in hashing with chaining to insert n elements into a table of size n the expected length of the longest chain is $\tilde{O}(n^{1/3})$. The proof also generalises to give the same bound when we use the multiply-shift hash function by Dietzfelbinger et al. [Journal of Algorithms 1997].

*Research partly supported by Advanced Grant DFF-0602-02499B from the Danish Council for Independent Research under the Sapere Aude research career programme and by the FNU project AlgoDisc - Discrete Mathematics, Algorithms, and Data Structures

1 Introduction

In this paper we study the hash function $h : [p] \rightarrow [m]$ (where $[m] = \{0, 1, \dots, m-1\}$) defined by $h(x) = ((ax + b) \bmod p) \bmod m$, where $a, b \in [p]$ are chosen uniformly at random from $[p]$. Here, p is a prime and $p \geq m$. We assume that we have a set $X \subseteq [p]$ of n keys with $n \leq m$ and use h to assign a hash value $h(x)$ to each key $x \in X$. We are interested in the frequency of the most popular hash value, i.e. we study the random variable $M(h, X)$ defined by

$$M(h, X) = \max_{y \in [m]} |\{x \in X \mid h(x) = y\}|. \quad (1)$$

In Theorem 1 we prove that $\mathbb{E}[M(h, X)] = O(\sqrt[3]{n \log n})$. We also consider the hash function $\bar{h} : [q] \rightarrow [m]$ defined by $\bar{h}(x) = \left\lfloor \frac{(ax) \bmod q}{q/m} \right\rfloor$, where q, m are powers of 2, $q \geq m \geq n$ and a is chosen uniformly at random among the odd numbers from $[q]$. The function $\bar{h}(x)$ was first introduced by Dietzfelbinger et al. [3]. In Theorem 2 we prove that it also holds that $\mathbb{E}[M(\bar{h}, X)] = O(\sqrt[3]{n \log n})$.

We note that when we use $h(x) = ((ax + b) \bmod p) \bmod m$ in hashing with chaining, M is the size of the largest chain. When scanning the hash table for an element the expected time used is $O(1)$ and the worst case time is at most $O(M(h, X))$.

1.1 Related work

It is folklore that the size of the largest chain is $O(\sqrt{n})$ and this bounds hold for any 2-independent hash function.

Alon et al. [1] considers the linear hash function $h_{m,k} : \mathcal{F}^m \rightarrow \mathcal{F}^k$, where \mathcal{F} is a finite field and $n = |\mathcal{F}|^k$. The function is defined by $h_{m,k}(x_1, \dots, x_m) = \sum_i x_i a_i$, where $a_i \in \mathcal{F}^k$ is chosen uniformly at random. For $m = 2, k = 1$ the hash function is $h_{2,1}(x, y) = ax + by$ where $a, b \in \mathcal{F}$ are chosen uniformly at random. It is shown in [1] that there exists a set $X \subseteq \mathcal{F}^2$ such that $\mathbb{E}[M(h_{2,1}, X)] > \sqrt{n}$ if n is a square and $\mathbb{E}[M(h_{2,1}, X)] = \Omega(\sqrt[3]{n})$ if n is a prime power that is not a square. In [1] it is also shown that when \mathcal{F} is the field of two elements the expected length of the longest chain is $O(\log n \log \log n)$ improving the results in [4, 5].

Broder et al. [2] considered $h(x) = (ax + b) \bmod p$ in the context of min-wise hashing.

2 Preliminaries

\mathbb{Z} denotes the integers, and $\mathbb{Z}_n = \mathbb{Z}/n\mathbb{Z}$ denotes the integers mod n . \mathbb{Z}_n^* is the set of elements of \mathbb{Z}_n having a multiplicative inverse. $[n]$ is the set of integers from 0 to $n-1$, that is $[n] = \{0, 1, 2, \dots, n-1\}$. For a pair of integers $n, m \in \mathbb{Z}$ such that $(n, m) \neq (0, 0)$ we let $\gcd(n, m)$ denote the greatest common divisor of n and m . If $\gcd(n, m) = 1$ then n and m are said to be *coprime*.

For integers $x, r \in \mathbb{Z}$ we let $[x]_r \in \mathbb{Z}_r$ denote the residue class of $x \bmod r$. We let $\iota_r : \mathbb{Z}_r \rightarrow [r]$ be the unique mapping that satisfies $[\iota_r(x)]_r = x$. For $x \in \mathbb{Z}_r$ we let $\|x\|_r = \min\{\iota_r(x), \iota_r(-x)\}$.

For a set S and an element x the sets $S + x$ and xS are defined as $\{s + x \mid s \in S\}$ and $\{xs \mid s \in S\}$, respectively.

For integers r, s, m , let $I_m(r, s) \subseteq \mathbb{Z}_m$ denote the set

$$I_m(r, s) = \{[r]_m, [r+1]_m, \dots, [r+s-1]_m\}.$$

The set $I_m(r, s)$ is called an *interval*. A non-empty set $X \subseteq \mathbb{Z}_m$ is an interval if there exists r, s such that $X = I_m(r, s)$.

3 Main Result

In this section we prove the main results of this paper, namely Theorem 1 and Theorem 2. The proofs of the two theorems are very similar and both rely on Lemma 1 below.

Lemma 1. *Let n, M, r be integers satisfying $4M \leq n \leq r$ and let $A \subseteq \mathbb{Z}_r$ be a set of size n . Let $B \subseteq \mathbb{Z}_p^*$ be a set of size $\leq M$ satisfying the following conditions:*

- 1 $\iota(b) \in (M, 2M)$ for all $b \in B$.
- 2 $\iota(b), \iota(b')$, and r are pairwise coprime for every $b, b' \in B$ with $b \neq b'$.

Assume that for every $b \in B$ there exists an interval I_b of size $\lceil \frac{r}{n} \rceil$ such that $I_b \cap bA$ contains at least $4M$ elements. Then there exists at least $M|B|$ ordered pairs of different elements $a, a' \in A$ such that $|a - a'| < \frac{r}{nM}$.

Proof. We note that for every b the set $b^{-1}I_b$ is the union of $\iota(b)$ disjoint intervals of size $\leq \lceil \frac{r}{n\iota(b)} \rceil$, and we write it as such a union $b^{-1}I_b = \bigcup_{j=0}^{\iota(b)-1} I_{b,j}$. For any $b, b' \in B, b \neq b'$ the set $b^{-1}I_b \cap b'^{-1}I_{b'}$ is either empty or an interval. So for each $b, b' \in B, b \neq b'$ there is at most one index $j \in [\iota(b)]$ such that the intersection $I_{b,j} \cap b'^{-1}I_{b'}$ is non-empty. For every $b \in B$ and $j \in [\iota(b)]$, let $\delta(b, j)$ denote the number of elements $b' \in B$ such that $I_{b,j} \cap b'^{-1}I_{b'}$ is non-empty. Note that $\delta(b, j) \geq 1$ since $b \in B$. Furthermore $\sum_{j=0}^{\iota(b)-1} \delta(b, j) < |B| + \iota(b) \leq 3M$ since each $b'^{-1}I_{b'}$ has a non-empty intersection with at most one of the sets $I_{b,j}, j \in [\iota(b)]$.

The number of ordered pairs of different elements $(a, a') \in A \cap I_{b,j}$ such that $|a - a'| < \frac{r}{nM}$ is exactly $|A \cap I_{b,j}| \cdot (|A \cap I_{b,j}| - 1)$ since $I_{b,j}$ is an interval of size $\leq \lceil \frac{r}{n\iota(b)} \rceil$ and $\iota(b) > M$. Let $\tau(b, j) = \max\{0, |A \cap I_{b,j}| - 1\}$, then the number of pairs is at least $(\tau(b, j))^2$. We can lower bound the number of such pairs in A by considering the pairs in $A \cap I_{b,j}$ for each $b \in B$ and $j \in [\iota(b)]$ and note that each pair we count is counted at most $\delta(b, j)$ times. This gives that the number of ordered pairs $(a, a') \in A$ such that $|a - a'| < \frac{r}{nM}$ is at least:

$$\sum_{b \in B} \sum_{j \in [\iota(b)]} \frac{(\tau(b, j))^2}{\delta(b, j)} \quad (2)$$

For any $b \in B$, by the Cauchy-Schwartz inequality we have that:

$$\left(\sum_{j \in [\iota(b)]} \delta(b, j) \right) \left(\sum_{j \in [\iota(b)]} \frac{(\tau(b, j))^2}{\delta(b, j)} \right) \geq \left(\sum_{j \in [\iota(b)]} \tau(b, j) \right)^2 \quad (3)$$

We clearly have that $\sum_{j \in [\iota(b)]} \tau(b, j) \geq 4M - \iota(b) \geq 2M$. Also recall, that we have that $\sum_{j \in [\iota(b)]} \delta(b, j) \leq 3M$. Combining this with (2) and (3) gives that A contains at least $\frac{4M|B|}{3} \geq M|B|$ of the desired pairs. \square

Below is a proof of Theorem 1.

Theorem 1. *Let n, m, p be integers with p a prime and $p \geq m \geq n$. Let $X \subseteq \mathbb{Z}_p$ be a set of n elements. Let $h : \mathbb{Z}_p \rightarrow \mathbb{Z}_m$ be defined by $h(x) = [\iota_p(ax + b)]_m$ where $a, b \in \mathbb{Z}_p$ are chosen uniformly at random. Let $M = M(X)$ be the random variable counting the number of elements $x \in X$ that hash to the most popular hash value, that is*

$$M = M(X) = \max_{y \in \mathbb{Z}_m} |\{x \in X \mid h(x) = y\}|.$$

Then

$$\mathbb{E}[M] = O\left(\sqrt[3]{n \log n}\right). \quad (4)$$

Proof. We note that $E[M \mid a = 0] = n$ since h is constant when $a = 0$. Therefore:

$$E[M] = \frac{p-1}{p} E[M \mid a \neq 0] + \frac{1}{p} E[M \mid a = 0] < E[M \mid a \neq 0] + 1.$$

Therefore it suffices to bound the expected value of M when a is chosen uniformly at random from $\mathbb{Z}_p \setminus \{0\} = \mathbb{Z}_p^*$ and not from \mathbb{Z}_p . So from now on, assume that a is chosen uniformly at random from \mathbb{Z}_p^* .

The random variables a and $a^{-1}b$ are independent. Note, that $h(x)$ can be rewritten as $h(x) = [\iota_p(a(x + a^{-1}b))]_m$. It clearly suffices to bound the expected value of M conditioned on all possible values $a^{-1}b$. For any fixed value of $a^{-1}b = c$, the expected value of M conditioned on $a^{-1}b = c$ is the same as the expected value of $M(X + c)$ conditioned on $b = 0$. Therefore it suffices to give the proof under the assumption that $b = 0$. So we assume that $b = 0$.

Let $A = m^{-1}aX$, then there exists an interval I_a of size at most $\lceil \frac{p}{m} \rceil$ that contains M elements of A for the following reason: Let $f : \mathbb{Z}_p \rightarrow \mathbb{Z}_m$ be defined by $x \rightarrow [\iota_p(x)]_m$. By definition, there exists a random variable $y \in \mathbb{Z}_p$ such that $|f^{-1}(y) \cap aX| \geq M$. And there exists a $i \in [m]$ such that

$$f^{-1}(y) = \left\{ [i + km]_p \mid k \in \mathbb{Z}, 0 \leq k < \frac{p-i}{m} \right\},$$

and hence $I_a = m^{-1}f^{-1}(y)$ is an interval of size $\leq \lceil \frac{p}{m} \rceil$ that contains M elements of A .

Let $\alpha \in [1, \frac{n}{4}]$. We are now going to bound the probability that $M \geq 4\alpha$. Let $\delta = \Pr[M \geq 4\alpha]$ and let \mathcal{A} be the set of all elements $a_0 \in \mathbb{Z}_p^*$ such that $M \geq 4\alpha$ if $a = a_0$.

Let $S \subseteq \mathbb{Z}_p^*$ be the set of all elements $s \in \mathbb{Z}_p^*$ that satisfies that $\iota_p(s)$ is a prime in the interval $(\alpha, 2\alpha)$. Let $B \subseteq S$ be the set of all elements $s \in S$ such that $as \in \mathcal{A}$. Note, that B is a random variable. By linearity of expectation, we have that $E[|B|] = |S| \delta$. Recall, that $A = m^{-1}aX$. For any $b \in B$ we have that $ab \in \mathcal{A}$ and therefore there exists an interval of size $\lceil \frac{x}{n} \rceil$ that contains at least 4α elements of bA . By Lemma 1, this implies that there are $\alpha |B|$ ordered pairs of different elements $x, x' \in X$ such that $\|ax - ax'\|_p < \frac{p}{m\alpha}$. So the expected number of elements $x, x' \in X$ such that $\|a(x - x')\|_p < \frac{p}{m\alpha}$ is at least $\alpha E[|B|] = \alpha \delta |S|$. On the other hand, for each ordered pair of different elements $x, x' \in X$ the probability that $\|a(x - x')\|_p < \frac{p}{m\alpha}$ is at most $\frac{2p}{m\alpha(p-1)}$, and by linearity of expectation the expected number of such ordered pairs is at most

$$n(n-1) \cdot \frac{2p}{m\alpha(p-1)} \leq \frac{2n}{\alpha}.$$

We conclude that $\alpha \delta |S| \leq \frac{2n}{\alpha}$. By the prime number theorem, $|S| = \Theta\left(\frac{\alpha}{\log \alpha}\right) = \Omega\left(\frac{\alpha}{\log n}\right)$. Reordering gives us that:

$$\Pr[M \geq 4\alpha] = \delta = O\left(\frac{n \log n}{\alpha^3}\right).$$

The expected value of M can now be bounded in the following manner:

$$\begin{aligned} E[M] &= \sum_{k=1}^{\infty} \Pr[M \geq k] \\ &= \sum_{k=1}^{\lfloor \sqrt[3]{n \log n} \rfloor} \Pr[M \geq k] + \sum_{k=\lfloor \sqrt[3]{n \log n} \rfloor + 1}^n \Pr[M \geq k] \\ &\leq \left\lfloor \sqrt[3]{n \log n} \right\rfloor + \sum_{k=\lfloor \sqrt[3]{n \log n} \rfloor + 1}^n O\left(\frac{n \log n}{k^3}\right) \\ &= O\left(\sqrt[3]{n \log n}\right) \end{aligned}$$

which was what we wanted. \square

The proof of Theorem 2 is very similar to the proof of Theorem 1 but we include it for completeness.

Theorem 2. *Let n, ℓ, r, q, m be integers with $q = 2^r, m = 2^\ell$ and $q \geq m \geq n$. Let $X \subseteq \mathbb{Z}_q$ be a set of n elements. Let $h : \mathbb{Z}_q \rightarrow [m]$ be defined by $h(x) = \lfloor \iota_q(ax) \cdot 2^{\ell-r} \rfloor$ where $a \in \mathbb{Z}_q^*$ are chosen uniformly at random. Let $M = M(X) = \max_{y \in [m]} |\{x \in X \mid h(x) = y\}|$. Then*

$$\mathbb{E}[M] = O\left(\sqrt[3]{n \log n}\right). \quad (5)$$

Proof. Let y be a random variable such that $|h^{-1}(y) \cap X| = M$, and let $A = aX$. The set $ah^{-1}(y)$ is an interval of size $\frac{q}{m}$ that contains exactly M elements of A .

Let $\alpha \in [1, \frac{n}{4}]$. We are now going to bound the probability that $M \geq 4\alpha$. Let $\delta = \Pr[M \geq 4\alpha]$, and let \mathcal{A} be the set of all elements $a_0 \in \mathbb{Z}_p^*$ such that $M \geq 4\alpha$ if $a = a_0$.

Let $S \subseteq \mathbb{Z}_q^*$ be the set of all elements $s \in \mathbb{Z}_q^*$ that satisfies that $\iota_q(s)$ is a prime in the interval $(\alpha, 2\alpha)$. Let $B \subseteq S$ be the set of all elements $s \in S$ such that $as \in \mathcal{A}$. Note, that B is a random variable. By linearity of expectation, we have that $\mathbb{E}[|B|] = |S| \delta$. Recall, that $A = m^{-1}aX$. For any $b \in B$ we have that $ab \in \mathcal{A}$ and therefore there exists an interval of size $\frac{q}{m}$ that contains at least 4α elements of bA . By Lemma 1, this implies that there are $\alpha |B|$ ordered pairs of different elements $x, x' \in X$ such that $\|ax - ax'\|_q < \frac{q}{m\alpha}$. So the expected number of elements $x, x' \in X$ such that $\|a(x - x')\|_q < \frac{q}{m\alpha}$ is at least $\alpha \mathbb{E}[|B|] = \alpha \delta |S|$. On the other hand, for each ordered pair of different elements $x, x' \in X$ the probability that $\|a(x - x')\|_q < \frac{q}{m\alpha}$ is at most $\frac{4}{m\alpha}$, and by linearity of expectation the expected number of such ordered pairs is at most

$$n(n-1) \cdot \frac{4}{m\alpha} \leq \frac{4n}{\alpha}.$$

We conclude that $\alpha \delta |S| \leq \frac{4n}{\alpha}$, and now we can bound the expected value exactly as in Theorem 1. \square

References

- [1] Noga Alon, Martin Dietzfelbinger, Peter Bro Miltersen, Erez Petrank, and Gábor Tardos. Linear hash functions. *Journal of the ACM (JACM)*, 46(5):667–683, 1999.
- [2] Andrei Z. Broder, Moses Charikar, Alan M. Frieze, and Michael Mitzenmacher. Min-wise independent permutations. *J. Comput. Syst. Sci.*, 60(3):630–659, 2000.
- [3] Martin Dietzfelbinger, Torben Hagerup, Jyrki Katajainen, and Martti Penttonen. A reliable randomized algorithm for the closest-pair problem. *Journal of Algorithms*, 25(1):19–51, 1997.
- [4] George Markowsky, Larry Carter, and Mark N. Wegman. Analysis of a universal class of hash functions. In *Mathematical Foundations of Computer Science 1978, Proceedings, 7th Symposium, Zakopane, Poland, September 4-8, 1978*, pages 345–354, 1978.
- [5] Kurt Mehlhorn and Uzi Vishkin. Randomized and deterministic simulations of prams by parallel machines with restricted granularity of parallel memories. *Acta Inf.*, 21:339–374, 1984.