

Appendix: Entropy Learned Hashing

10x Faster Hashing with Controllable Uniformity

BRIAN HENTSCHEL, Harvard University
UTKU SIRIN, Harvard University
STRATOS IDREOS, Harvard University

This document accompanies Entropy-Learned Hashing and covers aspects of the paper which were not covered in the main body of the paper. This covers proofs regarding Entropy-Learned Hashing and linear probing, discussion of robustness for Entropy-Learned Hashing data structures, and additional experiments.

APPENDIX A PROOFS FOR LINEAR PROBING

A.1 Notation & Problem Setup

The notation and problem setup are mostly the same as in the original paper. We give a brief recap of the notation and problem setup here for convenience.

Hash functions in Entropy-Learned hashing consists of two parts, L and H . The first part L chooses which locations of an input item x to hash and maps a key x to any subkey of x . The second part is a general hash function H which gets applied to $L(x)$ so that the full partial key hash $H' = H \circ L$. It is assumed that H follows the model of a perfectly random hash function, so that each unique input is given an independent uniform output between 0 and $m - 1$.

To analyze algorithms which use partial-key hashing, we need to deal with the fact that the keys entered into H are no longer unique. Towards this end we introduce notation for the set of keys K contained in a hash based data structure. The multi-set $S_{|L} = (K_{|L}, z)$ is made of $K_{|L}$, the set of all partial-keys (outputs of L applied to keys in K), and z , which maps each key in $K_{|L}$ to the cardinality of its pre-image in K . For instance, if L takes the first two characters of an input and $K = \{\text{dog}, \text{dot}, \text{cat}, \text{fan}\}$, then $K_{|L} = \{\text{"do"}, \text{"ca"}, \text{"fa"}\}$, $z(\text{"ca"}) = 1$, and $z(\text{"do"}) = 2$. We use as shorthand throughout z_x for $z(x)$.

Finally, we assume some familiarity with the linear probing algorithm, wherein collisions are resolved by checking the next slot in a circular array until an empty one is found. In this document, we resolve collisions by going to the left, i.e. checking $i, i - 1, i - 2, \dots$ rather than $i, i + 1, \dots$. This is to be consistent with Knuth's proof for linear probing given in [3]. A reminder of the rest of the notation can be seen in the table below.

A.2 Linear Probing

As in the main document, we start our proof by going over full-key hashing, then analyze partial-key hashing with a fixed dataset, and finally cover partial-key hashing with random data. Before starting the proof, we briefly cover the results. When the multi-set $S_{|L}$ is fixed, the expectation is an expectation over the randomness of the hash function H . When it is not fixed, the expectation is taken over both the random multi-set $S_{|L}$ and the random hash function H .

A.2.1 Results For full-key hashing, given any set of n keys and a hash table with m slots, the expected number of comparisons for querying a key not in the dataset and the average number of comparisons when querying for

Authors' addresses: Brian Hentschel, Harvard University, bhentschel@g.harvard.edu; Utku Sirin, Harvard University, wilsonqin@seas.harvard.edu; Stratos Idreos, Harvard University, stratos@seas.harvard.edu.

Notation	Definition
X, x	keys stored in the filter or hash table
Y, y	key which is queried for
H	hash function
m	size of table (in slots)
n	number of keys in table
K	set of keys
$S_{ L}$	multi-set of keys conditioned on L . Equal to $(K_{ L}, z)$ where $K_{ L}$ consists of all partial keys, and z maps each key $x \in K_{ L}$ to the number of times it appears in K . z_x will be used as shorthand for $z(x)$ throughout.
α	fill of hash table: $\frac{n}{m}$
P'	number of comparisons to find non-existing key
P	average # of comparisons to retrieve an existing key

Table 1. Notation used throughout the paper

an existing key are:

$$\mathbb{E}[P'] \leq \frac{1}{2} \left(1 + \frac{1}{(1-\alpha)^2} \right) \quad (1)$$

$$\mathbb{E}[P] \leq \frac{1}{2} \left(1 + \frac{1}{1-\alpha} \right) \quad (2)$$

Letting $c = \sum_x z_x^2$ be the number of collisions in $S_{|L}$ and $d = \sum_{x: z_x \geq 2} z_x$ the number of duplicated keys in $S_{|L}$, the costs for partial-key hashing are

$$\mathbb{E}[P'] \leq \begin{cases} \frac{1}{2} \left(1 + \frac{1}{(1-\alpha)^2} + \sum_{x \neq y} \frac{z_x^2}{m(1-\alpha)^2} \right) & \text{if } z_y = 0 \\ \frac{z_y}{1-\alpha} + \frac{1}{(1-\alpha)^2} + \sum_{x \neq y} \frac{z_x^2}{m(1-\alpha)^2} & \text{if } z_y > 0 \end{cases} \quad (3)$$

$$\begin{aligned} \mathbb{E}[P] &\leq \frac{n-d}{2n} + \frac{1}{2} Q_0(m, n) + \frac{c}{m} Q_0(m, n) + \frac{c+d}{2n} Q_0(m, d) \\ &\approx \left(\frac{1}{2} + \frac{c}{n} \right) \left(1 + \frac{1}{1-\alpha} \right) \end{aligned} \quad (4)$$

where here y is the queried for key in P' and the approximation in (4) uses the assumption that d is small compared to m . For random data from an i.i.d. source with Renyi entropy H_2 , these translate to the following bounds:

$$\begin{aligned} \mathbb{E}[P'] &\leq \frac{1}{2} \left(1 + \frac{1}{(1-\alpha)^2} \right) + n 2^{-H_2(L(X))} \frac{3}{2(1-\alpha)^2} \\ \mathbb{E}[P] &\leq \frac{1}{2} \left(1 + \frac{1}{1-\alpha} \right) + n 2^{-H_2(L(X))} \left(1 + \frac{1}{1-\alpha} \right) \end{aligned}$$

Comparing Partial-Key Hashing with Full-Key Hashing. The above equations, when fully analyzed, make a rather intuitive point: when there are only a few duplicates in K_L so that it closely approximates full-key hashing, the number of comparisons required by full-key and partial-key hashing are close. The general form of each equation for partial key hashing is that it is the same as full-key hashing plus a small penalty term. These penalty terms go to 0 as the keys in partial-key hashing become more distinct.

For all equations, we note the equations above are relatively tight bounds as long as α is not close to 1. This is reasonable to assume as all hash tables keep $\alpha < 0.9$ in practice because of the extremely bad behavior of linear

probing tables as $\alpha \rightarrow 1$, wherein they degrade to linear scans. The most common range for α is between 0.25 and 0.85, with higher alpha leading to better memory usage but slower query times.

A.2.2 Analysis of Full Key Linear Probing To prove (3) and (4), we start by proving (1) and (2), first proven by Donald Knuth. Our proof initially follows steps taken in The Art of Computer Programming[3], however we deviate from the proof given there because that approach does not generalize to partial-key hashing. Additionally, we believe this proof is simpler to follow than the proof in [3].

Counting Hash Sequences. For n distinct items, we have m^n possible ways to assign them to $[m] = \{0, \dots, m-1\}$, all of which are equally likely by our hashing assumptions. Of all possible hash sequences, we first consider how many leave position 0 empty in the hash table.

THEOREM 1. *The number of hash sequences such that n items are hashed into $[1, m-1]$, and no overflow occurs so that position 0 is empty is*

$$\begin{cases} (1 - \frac{n}{m}) \cdot h(m, n) & \text{if } 0 \leq n \leq m \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

where $h(m, n)$ is the number of ways to hash n objects into $[0, m-1]$.

This is a slight reformulation of the way it is stated in [3] so that it generalizes to partial-key hashing. To prove the theorem, note that hash sequences which hash n items into $[1, m-1]$ and do not overflow into position 0 are precisely the same as those that leave position 0 empty when hashing into a hash table of size m and resolving collisions via linear probing. This probability is $1 - \frac{n}{m}$ by the circular symmetry of linear probing. Thus, for full-key hashing there are $(1 - \frac{n}{m})m^n$ sequences which leave location 0 empty.

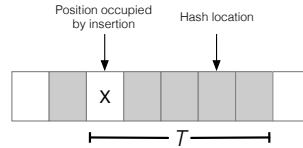
The next theorem gives the probability that a run of occupied slots begins at location 1 in the table and stops at location k .

THEOREM 2. *Let $g(m, n, k)$ be the number of sequences when hashing into a hash table of size m that leave position 0 empty, positions 1 to $k-1$ filled, and k empty. This is*

$$g(m, n, k) = \binom{n}{k-1} \frac{1}{k} \frac{m-n-1}{m-k} h(k, k-1) h(m-k, n-k+1) \quad (6)$$

This theorem simply builds off of Theorem 1 as there are $\binom{n}{k-1}$ ways to choose the $k-1$ elements in the run starting at 1, and then by Theorem 1 there are $\frac{1}{k} h(k, k-1)$ and $\frac{m-n-1}{m-k} h(m-k, n-k+1)$ ways to create the two subsequences with locations 0, k empty for each choice of the $k-1$ elements.

Analyzing Chain Length. We use the two theorems above to analyze a random variable T which represents the length counting from the empty position which a new item is inserted into to the final occupied position before the next empty position (going from left to right). The figure below shows an example.



If we know the expected value for T for a new item to be inserted, we can get its average distance from its hash location by the uniform randomness of hashing. Namely, for any chain of length T , each location in the chain is equally likely as an initial hash location, so $\mathbb{E}[P'|T] = \frac{1}{T} \sum_{i=1}^T i = \frac{1}{2} + \frac{1}{2}T$. It follows that

$$\mathbb{E}[P'] = \mathbb{E}[\mathbb{E}[P'|T]] = \frac{1}{2} + \frac{1}{2} \mathbb{E}[T]$$

To see how the two theorems above apply to our analysis of T , assume that the insertion key hashes to location a . If the insertion key is hashed into a chain of length t , then for some $t \geq 1, 0 \leq k \leq t-1$, we have $a-k$ empty, $a-k+1, \dots, a+(t-k-1)$ full, and $a+(t-k)$ empty. We note by symmetry that the same number of sequences produce a chain of this nature as which produce a chain such that 0 is empty, $1, \dots, t-1$ is full, t is empty. Since we have t choices for k , it then follows that there are $t \cdot g(m, n, t)$ hash sequences for which the insertion item is inserted into a chain of length t .

The approach taken in Knuth's art of computer programming is to analyze $E[T] = \sum_{t \geq 1} t^2 g(m, n, t)$ directly by using Abel's Binomial Theorem, i.e.

$$(x+y)^n = \sum_k \binom{n}{k} x(x-kz)^{k-1} (y+kz)^{n-k}$$

While correct, this approach feels magical as the use of Abel's Binomial is unintuitive (at least to the authors). Additionally, it does not generalize to the case of partial-key hashing, and so we take a different approach.

Analyzing Chain Length By Connecting Chain Length and the Probability of Chain Inclusion. Our approach to analyze $\mathbb{E}[T]$ is to look at the probability that each item is in the same probe chain as the inserted item. Namely, if we let x_1, \dots, x_n be the keys inserted into the hash table, then we have that

$$\mathbb{E}[T] = 1 + \sum_{i=1}^n P(x_i \in T)$$

Here we slightly abuse notation to have T both represent the length of the chain on the left and the set of objects in the chain on the right. We now connect this probability that $x_j \in T$, which is equal for all j , to the expected chain length conditioned on a single new item x being part of the chain T .

More concretely, let T_i be the chain length of an inserted item if i of the n items in the hash table have the same hash value as the item to be inserted. All other items are distributed uniformly at random and independently of the hash value of the to be inserted item. The following theorem holds.

THEOREM 3. *If C represents a set of i values which are fixed to have the same hash location as the newly inserted item, then for $x \notin C$, we have*

$$P(x \in T_i) = \frac{1}{m} \mathbb{E}[T_{i+1}]$$

PROOF. We first derive the formula for $\mathbb{E}[T_i]$. Assume as before that the insertion key hashes to location a . If the new item is hashed into a chain of length t , then for some $t \geq 1, 0 \leq k \leq t-1$, we have $a-k$ empty, $a-k+1, \dots, a+(t-k-1)$ full, and $a+(t-k)$ empty. Again by symmetry, the number of sequences that produce a chain of this nature is the same as the number of sequences that produce a chain such that 0 is empty, $1, \dots, t-1$ is full, t is empty, and the i items hash to location k . Let this number be $g(m, n, t, i, k)$. Then $f(m, n, t, i) = \sum_{k \geq 0} g(m, n, t, i, k)$ is the total number of sequences such that the new item is hashed into a chain of length t .

The value for $f(m, n, t, i)$ is calculable directly as this is identical to hashing $t-i-1$ items of multiplicity 1 and 1 value of multiplicity i into the first t slots keeping slot 0 empty, and hashing $n-(t-1)$ items into the final $m-t$ slots. Using theorem 1, there are $\frac{1}{t} h(t, t-i)$ ways to do the first and $\frac{m-n-1}{m-t} h(m-t, n+1-t)$ ways to do the second, and $\binom{n-i}{t-i-1}$ ways to which $t-i-1$ elements from $K \setminus C$ are the elements in the run starting at 1. It follows that

$$f(m, n, t, i) = \binom{n-i}{t-i-1} \frac{1}{t} \frac{m-n-1}{m-t} h(t, t-i) h(m-t, n+1-t)$$

The expected value of T_i is then $m^{-(n-i)} \sum_{t \geq 1} t \cdot f(m, n, t, i)$.

We now calculate the probability that $x \notin C$ is in T_i . To do so, we count the number of chains which contain x , and then divide by the total number of possible hash sequences. Using the same logic as before, for a key hashing to location a , the number of chains containing x is equivalent to the number of chains such that 0 is empty, $1, \dots, t-1$ is full and contains x , t is empty, and i items hash to location k . For a chain of length t , this means we have $\binom{n-i-1}{t-i-2}$ choices for the elements in the run between 0 and t , and we have $\frac{1}{t}h(t, t-i)$ and $\frac{m-n-1}{m-t}h(m-t, n-t+1)$ ways to hash these items to chain 1 and chain 2 respectively such that 0 and t are empty. Summing across all possible chain lengths, the number of chains which contain $x \notin C$ is

$$\begin{aligned} &= \sum_t \binom{n-i-1}{t-i-2} \frac{1}{t} \frac{m-n-1}{m-t} h(t, t-i) h(m-t, n-t+1) \\ &= \sum_t t f(m, n, t, i+1) \end{aligned}$$

where we use that $h(t, t-i-1) = t \cdot h(t, t-i)$. Dividing by the number of hash sequences, $m^{n-i} = m^{n-(i+1)}m$ gives that $P(x \in T_i) = \frac{1}{m} \mathbb{E}[T_{i+1}]$. \square

Finishing the Proof. Using these relationship between $\mathbb{E}[T_i]$ and $P(x \in T_i)$, we have

$$\begin{aligned} E[T] &= 1 + \sum_{i=1}^n P(x_i \in T_0) \\ &= 1 + \frac{n}{m} \mathbb{E}[T_1] \\ &= 1 + \frac{n}{m} (2 + \frac{n-1}{m} \mathbb{E}[T_2]) \end{aligned}$$

where here we use that $\mathbb{E}[T_i] = (i+1) + \sum_{x \notin C} P(x \in T_i)$. If we continue expanding the values of $\mathbb{E}[T_i]$, we get that

$$\begin{aligned} \mathbb{E}[T] &= 1 + 2\frac{n}{m} + 3\frac{n^2}{m^2} + \dots + \frac{(n+1)n!}{m^n} \\ &= Q_1(m, n) \end{aligned}$$

where

$$Q_r(m, n) = \sum_{k \geq 0} \binom{k+r}{k} \frac{n^k}{m^k}$$

It follows that

$$\mathbb{E}[P'] = \frac{1}{2} + \frac{1}{2} Q_1(m, n)$$

recovering the value proven in [3]. Noting that $\frac{n^k}{m^k} \leq \alpha^k$ and using that $\sum_{k \geq 0} (k+1)\alpha^k = d/d\alpha(\sum_{k \geq 0} \alpha^k) = (1-\alpha)^{-2}$ gives the desired bound $\mathbb{E}[P'] \leq \frac{1}{2}(1 + (1-\alpha)^{-2})$

Average Cost to Query an existing item. To go from $\mathbb{E}[P']$, the cost for a newly inserted key, to $\mathbb{E}[P]$, the average cost to query a key in the table, we average over the cost to insert the first n keys. Because the cost of an unsuccessful search in linear probing is the same as the cost to insert a key, this is the same as averaging

over unsuccessful searches over tables storing 1 through n elements. So $\mathbb{E}[P] = \frac{1}{n} \sum_{i \geq 0} \mathbb{E}[P'_i]$ where $\mathbb{E}[P'_i]$ is the expected cost to insert a key if there are i items in the table. Thus

$$\begin{aligned} \mathbb{E}[P] &= \frac{1}{2} + \frac{1}{2n} \sum_{i=0}^{n-1} Q_1(m, i) \\ &= \frac{1}{2} + \frac{1}{2n} \sum_{i=0}^{n-1} \sum_{k \geq 0} (k+1) \frac{i^k}{m^k} \\ &= \frac{1}{2} + \frac{1}{2n} \sum_{k \geq 0} \sum_{i=0}^{n-1} (k+1) \frac{i^k}{m^k} \end{aligned}$$

Now we use the rules of "finite calculus", which says that if $\Delta f(x) = f(x+1) - f(x) = g(x)$, then $\sum_{i=a}^b g(x) = f(b+1) - f(a)$. Here we have $\Delta x^k = kx^{k-1}$, so $\sum_{x=0}^{n-1} kx^{k-1} = n^k$. Plugging this in above, we have

$$\begin{aligned} \mathbb{E}[P] &= \frac{1}{2} + \frac{1}{2n} \sum_{k \geq 0} \frac{n^{k+1}}{m^k} \\ &= \frac{1}{2} + \frac{1}{2} Q_0(m, n-1) \end{aligned}$$

Using the same bound of $\frac{n^k}{m^k} \leq \alpha^k$, it follows that $E[P] \leq \frac{1}{2}(1 + \frac{1}{1-\alpha})$.

A.2.3 Proof of Costs Using Partial-Key Hashing To reason about the costs for partial-key hashing, we need to show how insertion and probe costs depend on S_L , the multi-set of keys. To clear up notation, we will refer to all multisets in theorems as C , C' , or C'' . For any multiset $C = (K, z)$, we have that $|C| = |K|$ is the number of distinct partial-key elements and $||C|| = \sum_{x \in C} z(x)$ is the total number of elements in C . We additionally extend our shorthand notation $z_x = z(x)$ to sets, i.e. $z_C = \sum_{x \in C} z_x = ||C||$. Additionally, we define $C_{2+} = x : z_x \geq 2$ to be the set of keys which are not unique in C .

Theorem 1,2 and 3 Restatements. Looking at theorem 1, it's proof holds for partial-key hashing. Namely, if C is a multi-set with $||C|| = n \leq m$, and C is hashed into locations $[1, m-1]$, then the number of hash sequences such that 0 is empty is

$$(1 - \frac{n}{m})h(m, C)$$

where $h(m, C)$ is the number of ways to hash the objects into $[0, m-1]$. For a given multiset C , we have $h(m, C) = m^{|C|}$.

To generalize theorem 2, we start by analyzing the form of equation (6) for $g(m, n, k)$. The equation consists of three parts: 1) ways to divide the set of n objects into one set with $k-1$ items and another with $n-k+1$ items, 2) ways to hash the $k-1$ and $n-k+1$ items into their arrays of size $k, m-k$, and 3) the scaling factors that say only $\frac{1}{k}, \frac{m-n-1}{m-k}$ of those leave slots 0 and k empty. It follows that when hashing a multi-set of items C , that

$$g(m, C, k) = \sum_{\substack{C' \subseteq C \\ ||C'||=k-1}} \frac{1}{k} \frac{m-n-1}{m-k} h(k, C') h(m-k, C \setminus C')$$

is the number of hash combinations that lead to a hash table with 0 empty, positions 1 to $k-1$ filled, and k empty.

Theorem 3 also holds for multisets, but we need new notation to be clearer about what items are included in the new chain. We define $T_{C'}$ to be the expected chain length of a new item if all items in the multiset $C' \subseteq C$ are guaranteed to have the same hash value as the newly inserted item.

THEOREM 4. *Let C be the multi-set of items being hashed into the hash table and let C' be a multi-set of items such that the items in C' are fixed to have the same hash value as a newly inserted item for the hash table. Let $x \notin C'$, and define $C'' = C' \cup \{(x, z_x)\}$. Then*

$$P(x \in T_{C'}) = \frac{1}{m} \mathbb{E}[T_{C''}]$$

PROOF. We start by deriving the equation for $\mathbb{E}[T_{C''}]$. Using the same reduction as in the proof of Theorem 3, the number of sequences for which the a new item is hashed into a chain of length t , given that the items in C'' hash to the same location as the new item, is the same as the number of chains with 0 empty, $1, \dots, t-1$ full, t empty, and for which the items in C'' hash to some location k between 0 and $t-1$. For a fixed chain $C^* \subset (C \setminus C'')$ with $\|C^*\| = t-1 - \|C''\|$, the number of ways to hash C^* and C'' into $0, \dots, t-1$ with the elements of C'' fixed to have the same hash location is equivalent to the number of ways to hash a new multi-set consisting of $C^* \cup (a, z_{C''})$ for an arbitrary new element a . By theorem 1 for multi-sets, it follows that the number of hash combinations such that $C^* \cup C''$ are between 0 and $t-1$ and location 0 is empty is $(1 - \frac{t-1}{t})(h(t, C^*) \cdot t)$. It follows that $f(m, C, t, C'')$, the number of hash combinations that produce a chain of length t , satisfies

$$f(m, C, t, C'') = \sum_{\substack{C^* \subset (C \setminus C'') \\ \|C^*\| = t-1 - \|C''\|}} \frac{m-n-1}{m-t} h(t, C^*) h(m-t, C \setminus (C^* \cup C''))$$

There are $m^{|C \setminus C''|}$ possible hash sequences and so it follows that $E[T_{C''}] = m^{-|C \setminus C''|} \sum_{t \geq 1} t f(m, C, t, C'')$.

We now derive the equation for $P(x \in T_{C'})$. The same logic again applies, with the exception that now x has a degree of freedom in its hash location. So for a fixed chain $C^* \subset (C \setminus C'')$ with $\|C^*\| = t-1 - \|C''\|$, the number of ways to hash C^* , x , and C' into $0, \dots, t-1$ with the elements of C' fixed to have the same hash location is equivalent to the number of ways to hash a new multi-set consisting of $C^* \cup \{(a, z_{C'}), (x, z_x)\}$ for an arbitrary new element a . Again by theorem 1 for multi-sets, this implies the number of hash combinations where x, C' , and C^* are hashed between 0 and $t-1$ with location 0 empty is $(1 - \frac{t-1}{t})(h(t, C^*) \cdot t^2)$. It follows that

$$\begin{aligned} P(x \in T_{C'}) &= m^{-|C \setminus C'|} \sum_{t \geq 1} \sum_{\substack{C^* \subset (C \setminus C'') \\ \|C^*\| = t-1 - \|C''\|}} \frac{m-n-1}{m-t} t \cdot h(t, C^*) h(m-k, S_{[L]} \setminus C^*) \\ &= m^{-1} m^{-|C \setminus C''|} \sum_{t \geq 1} t \cdot f(m, C, t, C'') \\ &= \frac{1}{m} E[T_{C''}] \end{aligned}$$

□

Bounds on $\mathbb{E}[T_{C'}]$. Theorem 4 again gives us a bridge between the probability of an item being in the same chain as a new item and the expected chain length when that item is hashed to the same location as a new item. We now use Theorem 4 to prove the following bound by induction on the set of unallocated objects, i.e. $C \setminus C'$:

$$\mathbb{E}[T_{C'}] \leq z_{C'} Q_0(m, n - z_{C'}) + Q_1(m, n - z_{C'}) + \sum_{x \in (C \setminus C')_{2+}} \frac{z_x^2}{m} Q_1(m, n - z_{C'})$$

Base case: Our base case is all sets such that $C' = C$, i.e. all objects in C are guaranteed to hash to the same location as the newly inserted item. The length of T is guaranteed to be $z_{C'} + 1 \leq z_{C'} Q_0(m, 0) + Q_1(m, 0)$.

Induction Steps: Assume our induction hypothesis holds for all sets $C, C' \subset C$ with $|C \setminus C'| \leq k$. Assume now that are given sets C, C' such that $|C \setminus C'| = k+1$, and that we wish to calculate $E[T_{C'}]$. As shorthand, let $C^- = C \setminus C'$. Then

$$\begin{aligned}
\mathbb{E}[T_{C'}] &\leq (1 + z_{C'}) + \sum_{x \in C^-} z_x P(x \in T_{C'}) \\
&= (1 + z_{C'}) + \frac{1}{m} \sum_{x \in C^-} z_x E(T_{C' \cup \{(x, z_x)\}}) \\
&\leq (1 + z_{C'}) + \frac{1}{m} \sum_{x \in C^-} z_x [(z_x + z_{C'}) Q_0(m, n - z_{C'} - z_x) + Q_1(m, n - z_{C'} - z_x) + \sum_{y \in C^-, y \neq x} \frac{z_y^2}{m} Q_1(m, n - z_{C'} - z_x)] \\
&\leq (1 + z_{C'}) + \frac{n - z_{C'}}{m} z_{C'} Q_0(m, n - z_{C'} - 1) + \sum_{x \in C^-} \frac{z_x^2 + z_x}{m} Q_0(m, n - z_{C'} - 1) \\
&\quad + \frac{n - z_{C'}}{m} Q_1(m, n - z_{C'} - 1) + \sum_{x \in C^-} \frac{z_x^2}{m} \frac{n - z_{C'}}{m} Q_1(m, n - z_{C'} - 1) \\
&= z_{C'} (1 + \frac{n - z_{C'}}{m} Q_0(m, n - z_{C'} - 1)) \tag{7} \\
&\quad + 1 + \frac{n - z_{C'}}{m} Q_0(m, n - z_c - 1) + \frac{n - z_{C'}}{m} Q_1(m, n - z_{C'} - 1) \tag{8} \\
&\quad + \sum_{x \in C_{2+}^-} \frac{z_x^2}{m} (Q_0(m, n - z_c - 1) + \frac{n - z_{C'}}{m} Q_1(m, n - z_{C'} - 1)) \tag{9}
\end{aligned}$$

Now, using the rules $\frac{n}{m} Q_1(m, n - 1) = Q_1(n, m) - Q_0(n, m)$ and $\frac{n}{m} Q_0(m, n - 1) = Q_0(m, n) - 1$ we can reduce equations (7), (8), and (9) respectively to $z_{C'} Q_0(m, n - z_{C'})$, $Q_1(m, n - z_{C'})$, and $\sum_{x \in (C \setminus C')_{2+}} \frac{z_x^2}{m} Q_1(m, n - z_{C'})$ respectively. This gives the required result on $\mathbb{E}[T_{C'}]$.

It follows for a query on a new item k such that $L(k) = y$, that it has a bound on its expected chain length of

$$E[T] \leq z_y Q_0(m, n - z_y) + Q_1(m, n - z_y) + \sum_{x \neq y \in C} \frac{z_x^2}{m} Q_1(m, n - z_y)$$

Connecting $\mathbb{E}[T]$ to $\mathbb{E}[P']$. Depending on whether $z_y = 0$, the connection between chain length and probe length changes. For items which match no other partial-keys in the dataset, we can use the same uniformity assumptions as before to recover that $\mathbb{E}[P'] = \frac{1}{2} + \frac{1}{2} \mathbb{E}[T]$. For items which match at least one partial-key, this uniformity assumption does not hold, and indeed items with larger c_y values are more likely to be hashed towards the beginning of their chain. Thus, we use that $P \leq T$ to say $\mathbb{E}[P] \leq \mathbb{E}[T]$, and note this is a loose bound when z_y is small. The result, overall, is equation 3, restated here for convenience.

$$\mathbb{E}[P'] \leq \begin{cases} \frac{1}{2} \left(1 + \frac{1}{(1-\alpha)^2} + \sum_{x \neq y} \frac{z_x^2}{m(1-\alpha)^2} \right) & \text{if } z_y = 0 \\ \frac{z_y}{1-\alpha} + \frac{1}{(1-\alpha)^2} + \sum_{x \neq y} \frac{z_x^2}{m(1-\alpha)^2} & \text{if } z_y > 0 \end{cases}$$

Average query cost for an existing key. The average cost to query a key in linear probing is unaffected by the order in which keys are inserted. Thus, we may assume any insertion order we desire in order to ease the analysis of the average number of probes for existing keys.

Because our analysis in the prior section created looser bounds for $z_y > 0$, namely in that there was no uniformity assumption on where items were in a chain, we add all duplicate keys first and then all non-duplicate

keys after. Expanding this out, we have

$$\begin{aligned}
 \mathbb{E}[P] &\leq \frac{1}{n} \left(\sum_{x \in C_{2+}} \frac{z_x^2}{2} Q_0(m, d) + \sum_{i=0}^{d-1} [Q_1(m, i) + \sum_{x \in C_{2+}} \frac{z_x^2}{m} Q_1(m, i)] \right) + \frac{1}{2n} \left[\sum_{i=d}^{n-1} (1 + Q_1(m, i) + \sum_{x \in C_{2+}} \frac{z_x^2}{m} Q_1(m, i)) \right] \\
 &\leq \frac{n-d}{n} + \frac{1}{2} Q_0(m, n-1) + \frac{c+d}{2n} Q_0(m, d-1) + \frac{c}{2m} (Q_0(m, n-1) + \frac{d}{n} Q_0(m, d-1)) \\
 &\approx \frac{1}{2} \left(1 + \frac{1}{1-\alpha} \right) + \frac{c}{n} + \frac{c}{m} \frac{1}{1-\alpha} \\
 &\leq \left(\frac{1}{2} + \frac{c}{n} \right) \left(1 + \frac{1}{1-\alpha} \right)
 \end{aligned}$$

Here we use that $\alpha_d = \frac{d}{m} \approx 0$ so that $\frac{1}{1-\alpha_d} \approx 1$. Because we are interested in the case that duplicate items are rare, this is a very good approximation.

Random Data. To go from fixed data to random data, we condition via Adam's law. For querying for a missing key, we first rewrite the expectation as

$$\begin{aligned}
 \mathbb{E}[P' | y, S_L] &\leq \frac{1}{2} + \mathbb{1}_{z_y \neq 0} \left(\frac{z_y}{1-\alpha} - \frac{1}{2} \right) + \left(\frac{1}{2} + \frac{1}{2} \mathbb{1}_{z_y \neq 0} \right) \frac{1}{(1-\alpha)^2} + \left(\frac{1}{2} + \frac{1}{2} \mathbb{1}_{z_y \neq 0} \right) \sum_x \frac{z_x^2}{m} \frac{1}{(1-\alpha)^2} \\
 &\leq \frac{1}{2} \left(1 - \frac{1}{(1-\alpha)^2} \right) + \frac{\mathbb{1}_{z_y \neq 0} z_y}{1-\alpha} + \frac{\mathbb{1}_{z_y \neq 0}}{(1-\alpha)^2} + \sum_x \frac{z_x^2}{m} \frac{1}{(1-\alpha)^2}
 \end{aligned}$$

Then using that $E[\mathbb{1}_{z_y \neq 0}] \leq \mathbb{E}[z_y] = \mathbb{E}[\mathbb{1}_{z_y \neq 0} z_y]$, that $\mathbb{E}[z_y] \leq n2^{-H_2(L(X))}$ by the union bound, and that $E[\sum_{x \in S_L} z_x^2] = n^2 2^{-H_2}$, we have

$$\begin{aligned}
 \mathbb{E}[P'] &\leq \frac{1}{2} \left(1 + \frac{1}{(1-\alpha)^2} \right) + \frac{n2^{-H_2(L(X))}}{1-\alpha} + \frac{n2^{-H_2(L(X))}}{2(1-\alpha)^2} + \sum_x \frac{\alpha n 2^{-H_2(L(X))}}{(1-\alpha)^2} \\
 &\leq \frac{1}{2} \left(1 + \frac{1}{(1-\alpha)^2} \right) + n2^{-H_2(L(X))} \frac{3}{2(1-\alpha)^2}
 \end{aligned}$$

More straightforwardly,

$$\mathbb{E}[P] \leq \frac{1}{2} \left(1 + \frac{1}{1-\alpha} \right) + n2^{-H_2(L(X))} \left(1 + \frac{1}{1-\alpha} \right)$$

This concludes the analysis of linear probing, proving the initial equations given in the paper.

APPENDIX B ROBUSTNESS OF ENTROPY-LEARNED HASHING DATA STRUCTURES

Robustness. For most algorithms and data structures, including those listed here, there is a trade-off between their robustness and their expected performance. How to choose between the expected performance and robustness is highly dependent on the application; certain applications might accept a 20% better average performance for a 1% chance at 2X worse performance whereas others may not. For Entropy-Learned Hashing and other techniques which exist as a base layer of potential use to many applications, the goal is to make this trade-off minimal, so that large gains in expected performance come at little cost for robustness. For Entropy-Learned Hashing, the amount of robustness for the hashing task depends on 1) the assumptions of Entropy-Learned Hashing, and 2) the task at hand. We start by discussing the assumptions of Entropy-Learned Hashing and then discuss Entropy-Learned Hashing for hash tables, Bloom Filters, and partitioning individually.

Assumptions of Entropy-Learned Hashing. Entropy-Learned Hashing makes weak assumptions, namely that data which are somewhat random remain somewhat random. Thus, unlike the usual case for learning

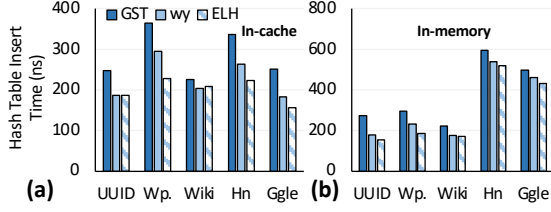


Fig. 1. Entropy-Learned Hashing reduces insert times for hash tables across datasets, and data sizes.

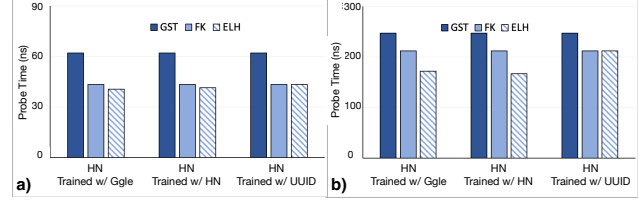


Fig. 2. Entropy-Learned Hashing reduces insert times for hash tables across datasets, and data sizes.

where shifts in data distributions almost always cause degradations in performance, shifts in data distribution for Entropy-Learned Hashing cause no change in performance as long as data stays random. In particular, if a distribution D_1 shifts to distribution D_2 this only causes a performance shift if $H_2(L(D_2)) < H_2(L(D_1))$. Thus the main robustness of Entropy-Learned Hashing comes from the fact that changes in distribution are unlikely to degrade performance. Still, we must address the case where D_2 actually becomes predictable on the bytes we choose to hash. We now cover this for each data structure.

Hash Tables. Entropy-Learned Hashing for hash tables is the most robust. This is for multiple reasons. First, if collisions are as expected on keys in the dataset, queries for both keys in the data and not in the data return quickly (see equations (4) and (5) for example). Second, we can monitor collisions during insertions with little overhead by counting the displacement from the initial hash position. Thus, this together with reason 1 means we have guaranteed good performance at all times. Third, Entropy-Learned Hashing can rehash if collisions ever deviate what is expected. The simplest way to do this is to revert to a full-key hash function, although other approaches such as backup bytes to include are an interesting future direction. Thus performance for hash tables can be guaranteed to be at least as good as traditional hashing (at the expense of code complexity) while usually being substantially better.

Bloom Filters. Bloom filters are less robust than hash tables. Their first defense in terms of robustness is the weak assumptions, namely that different distributions work for Entropy-Learned Hashing as long as the distribution is still random. Their second defense is that we can check that the data matches the desired distribution, namely because the # of set bits concentrates sharply around their expected value [1]. We use this fact to estimate the number of unique items in the filter. If data items are not as expected, or if queries are substantially different than the inserted items leading to a larger FPR, the filter must be rebuilt.

Partitioning. For partitioning, the cost of overloaded bins depends on the context, but for many contexts, such as in-memory radix partitioning, this can be solved by dividing the one or two overloaded bins into multiple bins. This is simple to implement and does not affect in-memory tasks like radix-partitioning much. For more complex tasks for which partitioning happens across a more expensive medium such as a network, approaches such as using the least loaded of d-bins would be of use [2].

APPENDIX C ADDITIONAL EXPERIMENTS

This Section covers the additional experiments described in Section 6.6. This includes experiments on 1) the efficiency of creating Entropy-Learned Hash data structures, 2) probing separate chaining hash tables, 3) experiments showing robustness properties, 4) experiments with dependent accesses (i.e. hash table lookups and Bloom filter lookups which must run one after the other instead of in parallel), and 5) additional experiments on Bloom filters showing different false positive rates.

1. Entropy-Learned Hashing Reduces Hash Table Insert Time. Entropy-Learned Hashing works similarly for inserts as it works for probes with hit rate = 1. Figure 1 presents hash table insert times for the five real-world

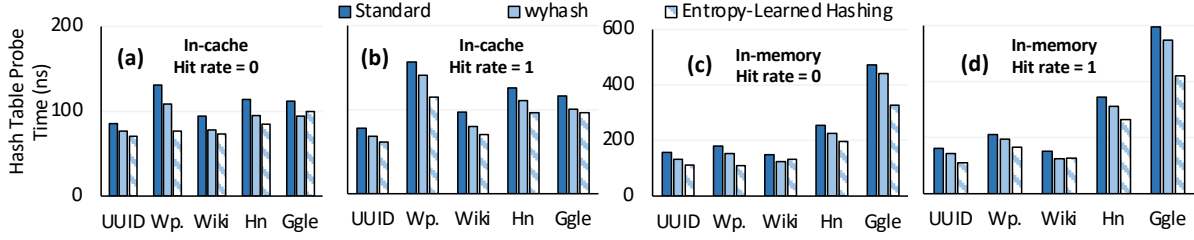


Fig. 3. Entropy-Learned Hashing reduces the probe time of standard chaining hash tables.

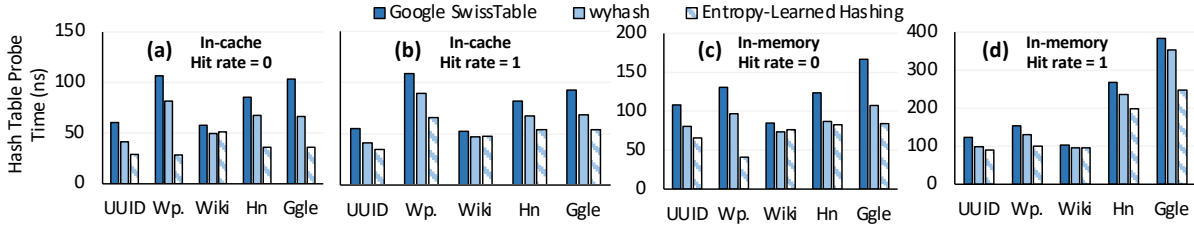


Fig. 4. Entropy-Learned Hashing reduces dependent probe times for hash tables across datasets, data sizes, and hit rates.

datasets we examine. The figure shows that Entropy-Learned hashing provides a $1.16\times$ to $1.3\times$ speedup over its closest competitor thanks to its reduced hash computation both for in-cache and in-memory data sizes. Thus, the speedups seen in query performance on Entropy-Learned hash tables carry over to performance on the insertion of data.

2. Separate Chaining Experiments. To show that Entropy-Learned Hashing works with separate chaining hash maps as well as linear probing hash maps, we integrate Entropy-Learned Hashing into `std::unordered_map`. The results can be seen in Figure 3. The key takeaway is that, as expected, Entropy-Learned Hashing reduces the probe times across datasets, data sizes, and hit rates, with this improvement being up to $1.72\times$. Thus Entropy-Learned Hashing works for separate chaining tables in addition to linear probing tables. The improvement is slightly lower than what is seen for SwissTable; this is mainly because `std::unordered_map` is much slower as a baseline and so portions of the hash table probe that are not the hash function take up more of the computation. A second change is that `std::unordered_map` is only a tiny amount faster when querying for non-existent keys than when querying for existent keys. This is again because `unordered_map` lacks optimizations present in SwissTable. Thus, a second takeaway of Figure 3 is that Entropy-Learned Hashing provides larger speedups for more optimized hash tables.

3. Robustness of Entropy-Learned Hashing. To test the robustness of Entropy-Learned hash tables, we test the performance of them when their training sets, i.e the set of sampled past queries and data items, does not match their test distribution of actual inserted and queried items. To test this, we vary the distribution given as a data sample and then insert and query items from the Hacker News dataset. Figures 2a,b show the results when using the large data size and querying for missing items (Fig. 2a) and existing items (Fig. 2b). As expected, when the Hacker News dataset is used to gather data as well as to insert and query items, we see speedups of 5% to 30%. More importantly, as Figures 2a and 2b show, when using Google URLs as a dataset to analyze ahead of time, using the Hacker News dataset to insert and query items still produces speedups of 5% to 27%. This shows the case where the test data is different in distribution but still random enough on the bytes of

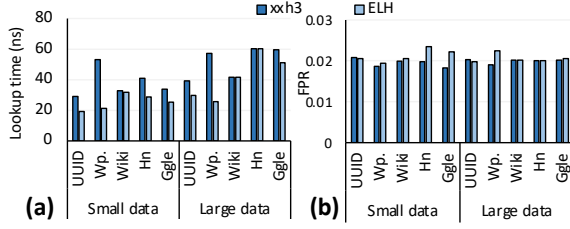


Fig. 5. Improving Bloom filter lookup time (a) and false positive rates (b) for small and large data sizes for dependent lookups.

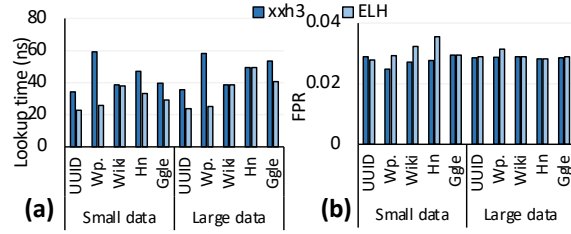


Fig. 6. Improving regular Bloom filter lookup time (a) and false positive rates (b) for small and large data sizes.

choice; in this case Entropy-Learned Hashing still provides speedups over full-key hashing. When using the UUID dataset as a training dataset, the user ids from UUID and posted urls from Hacker News are very different and Entropy-Learned Hashing defaults to using the full-key hash function as described in Section 2 of the technical report (and Section 5 of the main paper). Thus, in this case it provides no speedups but does not degrade performance.

4. Entropy-Learned Hashing Reduces Dependent Probe Times. Entropy-Learned Hashing reduces probe times for data structures when the probes are dependent (i.e. they cannot be batched). Figure 4 presents hash table probe times, and shows Entropy-Learned Hashing provides 1.18x to 2.9x speedups across different datasets, data sizes and hit rates. Figure 5 shows similar results for dependent Bloom filter lookups, where every lookup requires the result of the preceding lookup to start. The figure shows that Entropy-Learned Hashing significantly reduces the filter lookup times and provides a 1.16x to 2.5x speedup across datasets and data sizes.

The speedups are consistently lower for dependent lookups than they are for independent lookups. The reason is that independent lookups benefit from inter-lookup parallelism, as we discussed in Section 6.3. However, we observe a significant speedup even for dependent lookups. To illustrate, Entropy-Learned Hashing is 1.42x faster than wyhash for hash table probes on the Google dataset with large data and high hit rate. This is because there are two types of memory-level parallelism that Entropy-Learned Hashing benefits from: (i) inter-lookup parallelism, and (ii) intra-lookup parallelism. While independent hash table probes benefit from both inter- and intra-lookup parallelism, dependent hash table probes benefit from only intra-lookup parallelism. As a result, Entropy-Learned Hashing improves hash table probe time more when the probes are independent than it improves when the probes are dependent. To verify our hypothesis, we examined the MLP value for Entropy-Learned Hashing and wyhash for dependent hash table probes and found the MLP value is 2.0 for Entropy-Learned Hashing compared to 1.6 for wyhash. Thus ELH provides both parallelism benefits and computational benefits and produces consistent speedups for dependent lookups across small and large data sizes.

5. Experiments with different FPRs for Bloom Filters. The experiments in the main body of the text focused on throughput optimized filters, in particular using register-blocked Bloom filters from [4]. Figure 6 shows the performance of Entropy-Learned Hashing on traditional Bloom filters at a lower false positive rate of 1%. The results are extremely similar to Section 6.4, and so we don’t cover them in detail. Entropy-Learned Hashing again provides benefits in lookup performance of up to 2.4x at negligible (and tunable) increases to the false positive rate.

6. Entropy-Learned Hashing scales linearly with the number of threads. Up to now, we have focused on single-threaded behavior for Entropy-Learned Hashing. In this section, we examine how Entropy-Learned Hashing behaves as the number of threads is increased. We increase the number of threads concurrently probing

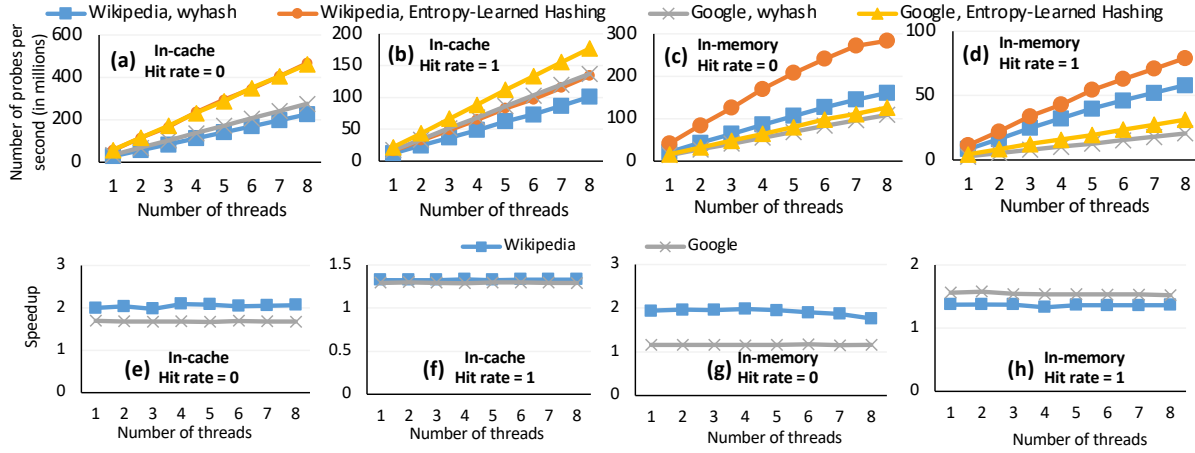


Fig. 7. Entropy-Learned Hashing scales linearly with the number of cores.

the same table and examine the number of probes per second by all the concurrently running threads. We pin each thread to a separate physical core without any hyper-thread sharing.

Figure 7 presents the results. The figure shows that both Entropy-Learned Hashing and wyhash scale linearly with the number of threads. As the number of threads is increased, the delivered throughput is increased in proportion with the number of threads. As a result, Entropy-Learned Hashing provides constant speedups over wyhash for all the number of threads.

REFERENCES

- [1] A. Broder, M. Mitzenmacher, and A. B. I. M. Mitzenmacher. Network applications of bloom filters: A survey. In *Internet Mathematics*, pages 636–646, 2002.
- [2] R. M. Karp, M. Luby, and F. Meyer auf der Heide. Efficient pram simulation on a distributed memory machine. In *Proceedings of the Twenty-Fourth Annual ACM Symposium on Theory of Computing*, STOC '92, page 318–326, New York, NY, USA, 1992. Association for Computing Machinery.
- [3] D. E. Knuth. *The Art of Computer Programming, Volume 3: (2nd Ed.) Sorting and Searching*. Addison Wesley Longman Publishing Co., Inc., USA, 1998.
- [4] H. Lang, T. Neumann, A. Kemper, and P. Boncz. Performance-optimal filtering: Bloom overtakes cuckoo at high throughput. *Proceedings of the VLDB Endowment*, 12(5):502–515, 2019.