

Lecture 26 — Hash Tables

Parallel and Sequential Data Structures and Algorithms, 15-210 (Spring 2012)

Lectured by Margaret Reid-Miller — 24 April 2012

Today:

- Hashing
- Hash Tables

1 Hashing

hash: transitive verb¹

1. (a) to chop (as meat and potatoes) into small pieces
(b) confuse, muddle
2. ...

This is the definition of hash from which the computer term was derived. The idea of hashing as originally conceived was to take values and to chop and mix them to the point that the original values are muddled. The term as used in computer science dates back to the early 1950's.

More formally the idea of hashing is to approximate a random function $h : \alpha \rightarrow \beta$ from a source (or universe) set U of type α to a destination set of type β . Most often the source set is significantly larger than the destination set, so the function not only chops and mixes but also reduces. In fact the source set might have infinite size, such as all character strings, while the destination set always has finite size. Also the source set might consist of complicated elements, such as the set of all directed graphs, while the destination are typically the integers in some fixed range. Hash functions are therefore many to one functions.

Using an actual randomly selected function from the set of all functions from α to β is typically not practical due to the number of such functions and hence the size (number of bits) needed to represent such a function. Therefore in practice one uses some pseudorandom function.

Exercise 1. *How many hash functions are there that map from a source set of size n to the integers from 1 to m ? How many bits does it take to represent them? What if the source set consists of character strings up length up to 20. Assume there are 100 possible characters.*

Why is it useful to have random or pseudo random functions that map from some large set to a smaller set? Generally such functions might be used to hide data, to reorder data in a random order, or to spread data smoothly over a range.

[†]Lecture notes by Guy E Blelloch, Margaret Reid-Miller, and Kanat Tangwongsan.

¹Merriam Websters

You have already seen one application of hashing, treaps. In particular we suggested using a hash function to hash the keys to generate the “random” priorities. Here, what was important is that the ordering of the priorities is somehow random with respect to the keys. Our analysis assumed the priorities were truly random, but it can be shown that a limited form of randomness that arise out of relatively simple hash functions is sufficient.

Today we will consider another common application, hashing to implement hash tables. In hash tables one is given a set of keys $K \subset U$ that need to be map to a range of integers so they can be stored at those locations in an array. The goal, here, is to spread the keys out across the integers as well as possible to minimize the number of keys that collide in the array. You should not confuse the terms hash function and hash table. They are two separate ideas, the later that uses the prior.

Next lecture we will consider a few other applications of hashing that make use of the hiding data aspect of hashing or a combination of hiding, reordering, and distributing data aspects. In addition, we will delay discussing ways to find appropriate hash functions until next lecture.

For now we will assume that we have hash functions with the following idealized property *simple uniform hashing*: The hash function uniformly distributes the n keys over the range $[0 \dots, m - 1]$ and the hash value for any key is independent of the hash value for any other key.

2 Hash Tables

Hash tables are used when an application needs to maintain a dynamic set where the main operations are `insert`, `find` and `delete`. Hash table can implement the abstract data types `Set` and `Table`. Unlike binary search trees, which require the universe of keys has a total ordering, hash tables do not. A total ordering enables the additional ordering operations provided by the `Ordered Set` abstract data type.

The main issue with hash table is collisions, where two keys hash to the same location. Is it possible to avoid collisions? Not if we don't know the set of keys in advance. Since the size of the table T is much smaller than the universe of keys U , $|T| \ll |U|$, there must exist at least two keys that map to the same index, by the *Pigeonhole Principle*: If you put more than m items into m bins, then at least one bin contains more than one item. For a particular hash function, the subset of keys $K \subset U$ that we want to hash may or may not cause a collision even when the number of keys is much smaller than the size of the table. Therefore, for general purpose dynamic hash tables we have to assume that collisions will occur.

How likely is there at least one collision? This is the same question as the birthday paradox: When there a n or more people in a room, what is the chance that two people have the same birthday. It turns out that for a table of size 365 you need only 23 keys for a 50% chance of a collision, and as little as 60 keys for a 99% chance. More interesting, when hashing to m locations, you can expect a collision after only $\sqrt{\frac{1}{2}\pi m}$ insertions, and can expect every location in the table has an element after $\Theta(m \log m)$ insertions. The former is related to the *birthday paradox*, whereas the later is relate to the *coupon collector* problem.

There are several well-studied collision resolution strategies:

- **Separate chaining**: Store elements not in a table, but in linked lists (containers, bins) hanging

off the table.

- **Open addressing:** Put everything into the table, but not necessarily into cell $h(k)$.
- **The perfect hash:** When you know the keys in advance, construct hash functions that avoids collisions entirely.
- **Multiple-choice hashing and Cuckoo hashing:** Consider exactly two locations $h_1(k)$ and $h_2(k)$ only.

We will consider the first two in this lecture, and save the others for another lecture. In our discussion we will assume we have a set of n keys K that we want to store and a hash function $h : \text{key} \rightarrow [1, \dots, m]$ for some m .

2.1 Separate Chaining

In 15-122 you covered hash tables using separate chaining. As you may recall, the idea is to maintain an array of linked lists. All keys that hashes to the same location in the sequence are maintained in a linked list. When inserting a key k , it inserts it at the beginning of the linked list at location $h(k)$. But if the application may attempt to insert the same key multiple times, separate chaining needs to search down the list to check whether the key is already in the list, in which case it might just as well add the key to the end of the linked list. To find a key, simply search for key in the linked list at location $h(k)$. To delete a key, remove it from the linked list.

The costs of these operations is related to the average length of a chain, which is n/m when there n keys in a table with m chains. We call $\lambda = n/m$ the *load factor* of the table.

We consider two cases when that can occur when applying these operations: an unsuccessful search when the key is not in the table, and a successful search when the key is in the table. We assume that we can compute the $h(k)$ in $O(1)$ work.

Claim 2.1. *Assuming simple uniform hashing, an unsuccessful search takes expected $\Theta(1 + \lambda)$ work.*

Proof. The average length of a list is λ and an unsuccessful search needs to search the whole list. Including the cost of computing $h(k)$, the total work is $\Theta(1 + \lambda)$. \square

Claim 2.2. *Assuming simple uniform hashing, a successful search takes expected $\Theta(1 + \lambda)$ work.*

Proof. To simplify the analysis, we assume that keys are added at the end of the chain². The the cost of a successful search is the same as an unsuccessful search at the time the key was added to the table. When the table has i keys, an unsuccessful search is expected to be $(1 + i/m)$. Averaging over all keys we get

$$\frac{1}{n} \sum_{i=0}^{n-1} (1 + i/m) = 1 + (n-1)/2m = 1 + \lambda/2 - \lambda/2m = \Theta(1 + \lambda)$$

\square

²The average successful search time is the same whether new keys are added to the front of the end of the chain

That is, successful searches examine half the list on average and unsuccessful searches full list. If $n = O(m)$ then with simple uniform hashing, all operations have expected $O(1)$ work and span. Even more important, some chains can be long, $O(n)$ in the worst case, but it is extremely unlikely they will more than double λ . The advantage of separate chaining is that it is not particularly sensitive to the size of the table. If the number keys is more than anticipated, the cost of search becomes only somewhat worse. If the number is less, then only a small amount of space in the table is wasted and cost of search is faster.

2.2 Open Address Hash Tables

The next technique does not need any linked lists but instead stores every key directly in the array. Open address hashing using so called linear probing has an important practical advantage over separate chaining: it causes fewer cache misses since typically all locations that are checked are on the same cache line.

The basic idea of open addressing is to maintain an array that is some constant factor larger than the number of keys and to store all keys directly in this array. Every cell in the array is either empty or contains a key.

To decide to which cells to assign a key, open addressing uses an ordered sequence of locations in which the key can be stored. In particular let's assume we have a function $h(k, i)$ that returns the i^{th} location for key k . We refer to the sequence $\langle h(k, 0), h(k, 1), h(k, 2), \dots \rangle$ as the *probe sequence*. We will get back to how the probe sequence might be assigned, but let's first go through how these sequences are used. When inserting a key the basic idea is to try each of the locations in order until it finds a cell that is empty, and then insert the key at that location. Sequentially insert would look like

```

1 fun insert(T, k) =
2   let
3     fun insert'(T, k, i) =
4       case nth T h(k, i) of
5         NONE => update(h(k, i), k) T
6       | _ => insert'(T, k, i + 1)
7   in
8     insert'(T, k, 1)
9   end

```

For example suppose the hash table has the following keys:

	0	1	2	3	4	5	6	7
T =		B			E	A		F

Now if for a key D we had the probe sequence $\langle 1, 5, 3, \dots \rangle$, then we would find location 1 and 5 full (with B and E) and place D in location 3 giving:

	0	1	2	3	4	5	6	7
T =		B		D	E	A		F

Note that, in order for the update operation to be constant work and span, T must be a single threaded array. Also, the code will loop forever if all locations are full. Such an infinite loop can be

prevented by ensuring that $h(k, i)$ tries every location as i is incremented, and checking when the table is full. Also, as given the code will insert the same key multiple times over. This problem is easily corrected by checking if the key is in the table and if so returning immediately.

To search we have the following code:

```

1  fun find(T,k) =
2  let
3    fun find'(T,k,i) =
4      case T[h(k,i)] of
5        NONE ⇒ false
6        | SOME(k') ⇒ if (eq(k,k')) then true
7                      else find'(T,k,i+1)
8  in
9    find'(T,k,1)
10 end

```

For example in the table above, if key E has the probe sequence $\langle 7, 4, 2, \dots \rangle$, `find` would first search location 7, which is full, and then location 4 to find E .

What if we want to delete a key? Let's say we deleted A from the table above, and then searched for D . Will `find` locate it? No, it will stop looking once it finds the empty cell where A was. One solution might be to rehash everything after a delete. But that would be an extremely expensive operation for every delete. An alternative is to use what is called a *lazy delete*. Instead of deleting the key, simply replace the key with a special HOLD value. That is, introduce an entry data type:

```

1  datatype  $\alpha$  entry = EMPTY | HOLD | FULL of  $\alpha$ 

```

Of course, we will need to change `find` and `insert` to handle HOLD entries. For `find`, simply skip over a HOLD entry and move to the next probe. Whereas, `insert(v)` can replace the first HOLD entry with `FULL(v)`. But if `insert` needs to check for duplicate keys, it first needs to search for the key. If it finds the key it overwrites it with the new value. Otherwise it continues until it finds an empty cell, at which point it can replace the first HOLD in the probe sequence.

The main concern with lazy deletes is that they effectively increases the load factor, increasing the cost of the hash table operations. If the load factor becomes large and performance degrades, the solution is to rehash everything to a new larger table. The table should be a constant fraction larger each time the table grow so as to maintain amortized constant costs.

Now let's consider some possible probe sequences we can use. Ideally, we would like a key to use any of the possible $m!$ probe sequences with equal probability. This ideal is called *uniform hashing*. But uniform hashing is not practical. Common probe sequences, which we will consider next, are

- linear probing
- quadratic probing
- double hashing

2.2.1 Linear Probing

In linear probing, to insert a key k , it first checks $h(k)$ and then checks each following location in succession, wrapping around as necessary, until it finds an empty location. That is, the i^{th} probe is

$$h(k, i) = (h(k) + i) \bmod m.$$

Each position determines a single probe sequence, so there are only m possible probe sequences.

The problem with linear probing is that keys tend to cluster. It suffers from *primary clustering*: Any key that hashes to any position in a cluster (not just collisions), must probe beyond the cluster and adds to the cluster size. Worse yet, primary clustering not only makes the probe sequence longer, it also makes it more likely that it will be lengthen further.

What is the impact of cluster for unsuccessful search? Let's consider two extreme examples when the table is half full, $\lambda = 1/2$. Clustering is minimized when every other location in the table is empty. In this case, the average number of probes needed to insert a new key k is $3/2$: One probe to check cell $h(k)$, and with probability $1/2$ that cell is full and it needs to look at the next location which, by construction, must be empty. In the worst case, all the keys are clustered, let's say at the end of the table. If k hashes to any of the first n locations, only one probe is needed. But hashing to the n^{th} location would require probing all n full locations before finally wrapping around to find an empty location. Similarly, hashing to the second full cell, requires probing $(n - 1)$ full cells plus the first empty cell, and so forth. Thus, under uniform hashing the average number of probes needed to insert a key would be

$$1 + [n + (n - 1) + (n - 2) + \dots + 1]/m = 1 + n(n + 1)/2m \approx n/4$$

Even though the average cluster length is 2, the cost for an unsuccessful search is $n/4$. In general, each cluster j of length n_j contributes $n_j(n_j + 1)/2$ towards the total number of probes for all keys. Its contribution to the average is proportional the *squares* of the length of the cluster, making long cluster costly.

We won't attempt to analyze the cost of successful and unsuccessful searches, as considering cluster formation during linear probing is quite difficult. We make the following claim:

Claim 2.3. When using linear probing in a hash table of size m that contains $n = \lambda m$ keys, the average number of probes needed for an unsuccessful search or an insert is

$$\frac{1}{2} \left(1 + \frac{1}{1 - \lambda^2} \right)$$

and for a successful search is

$$\frac{1}{2} \left(1 + \frac{1}{1 - \lambda} \right).$$

As you can see from the following table, which show the expected number of probes under uniform hashing, linear probing degrades significantly when the load factor increases:

λ	1/4	1/2	2/3	3/4	9/10
successful	1.2	1.5	2.0	3.0	5.5
unsuccessful	1.4	2.5	5.0	8.5	55.5

Linear probing is quite competitive, though, when the load factors are in the range (30-70%) as clusters tend to stay small. In addition, a few extra probes is mitigated when sequential access is much faster than random access, as is the case of caching. Because of primary clustering, though, it is sensitive to quality of the hash function or the particular mix of keys that result in many collisions or clumping. Therefore, it may not be a good choice for a general purpose hash tables.

2.2.2 Quadratic Probing

Quadratic probe sequences causes probes to move away from cluster, by making increasing larger jumps. The i^{th} probe is

$$h(k, i) = (h(k) + i^2) \bmod m.$$

Although, quadratic probing voids primary clustering, it still has *secondary clustering*: When two keys hash to the same location, they have the same probe sequence. Since there are only m locations in the table, there are only m possible probe sequences.

One problem with quadratic probing is that probe sequences do not probe all locations in the table. But since there are $(p + 1)/2$ quadratic residues when p is prime, we can make the following guarantee.

Claim 2.4. : *If m is prime and the table is at least half empty, then quadratic probing will always find an empty location. Furthermore, no locations are checked twice.*

Proof. (by contradiction) Consider two probe locations $h(k) + i^2$ and $h(k) + j^2, 0 \leq i, j < \lceil m/2 \rceil$. Suppose the locations are the same but $i \neq j$. Then

$$\begin{aligned} h(k) + i^2 &\equiv (h(k) + j^2) \bmod m \\ i^2 &\equiv j^2 \bmod m \\ i^2 - j^2 &\equiv 0 \bmod m \\ (i - j)(i + j) &\equiv 0 \bmod m \end{aligned}$$

Therefore, either $i - j$ or $i + j$ are divisible by m . But since both $i - j$ and $i + j$ are less than m and m is prime, they cannot be divisible by m . Contradiction.

Thus the first $\lceil m/2 \rceil$ probes are distinct and guaranteed to find an empty location. \square

Computing the next probe is only slightly more expensive than linear probing as it can be computed without using multiplication:

$$\begin{aligned} h_i - h_{i-1} &\equiv (i^2 - (i-1)^2) \bmod m \\ h_i &\equiv (h_{i-1} + 2i - 1) \bmod m \end{aligned}$$

Unfortunately, requiring that the table remains less than half full makes quadratic probing space inefficient.

2.2.3 Double Hashing

Double hashing uses two hash functions, one to find the initial location to place the key and a second to determine the size of the jumps in the probe sequence. The i^{th} probe is

$$h(k, i) = (h_1(k) + i \cdot h_2(k)) \bmod m.$$

Keys that hash to the same location, are likely to hash to a different jump size, and so will have different probe sequences. Thus, double hashing avoids secondary clustering by providing as many as m^2 probe sequences.

How do we ensure every location is checked? Since each successive probe is offset by $h_2(k)$, every cell is probed if $h_2(k)$ is relatively prime to m . Two possible ways to ensure $h_2(k)$ is relatively prime to m are, either make $m = 2^k$ and design $h_2(k)$ so it is always odd, or make m prime and ensure $h_2(k) < m$. Of course, $h_2(k)$ cannot equal zero.

Double hashing behaves quite closely to uniform hashing for careful choices of h_1 and h_2 . Under uniform hashing the average number of probes for an unsuccessful search or an insert is at most

$$1 + \lambda + \lambda^2 + \dots = \left(\frac{1}{1 - \lambda} \right)$$

and for a successful search is at most

$$\frac{1}{\lambda} \left(1 + \ln \left(\frac{1}{1 - \lambda} \right) \right).$$

The former bound is because the probability of needing more than i probes is at most λ^i . A search always needs one probe, and with probability λ needs a second probe, and with probability λ^2 needs a third probe, and so on. The bound for a successful search for a key k follows the same probe sequences as when it was first inserted. So if k was the $(j + 1)^{th}$ key inserted the cost of the for inserting it is at most $1/(1 - j/m)$. Therefore the average cost of a successful search is at most

$$\begin{aligned} \frac{1}{n} \sum_{j=0}^{n-1} \frac{1}{1 - j/m} &= \frac{m}{n} \sum_{j=0}^{n-1} \frac{1}{m - j} \\ &= \frac{1}{\lambda} (H_m - H_{m-n}) \\ &\leq \frac{1}{\lambda} (\ln m + 1 - \ln(m - n)) \\ &= \frac{1}{\lambda} \left(1 + \ln \left(\frac{1}{1 - \lambda} \right) \right) \end{aligned}$$

The table below shows the expected number of probes under the assumption of uniform hashing and is the best one can expect by open addressing.

λ	1/4	1/2	2/3	3/4	9/10
successful	1.2	1.4	1.6	1.8	2.6
unsuccessful	1.3	1.5	2.0	3.0	5.5

Comparing these numbers with the numbers in the table for linear probing, the linear probing numbers are remarkable close when the load factor is 50% or below.

The main advantage with double probing is that it allows for smaller tables (higher load factors) than linear or quadratic probing, but at the expense of higher costs to compute the next probe. The higher cost of computing the next probe may be preferable than longer probe sequences, especially when testing two keys equal is expensive.

2.3 Hash Table Summary

Hashing is a classic example of a space-time tradeoff: increase the space so table operations are faster, decrease the space but table operations are slower.

Separate chaining is simple to implement and is less sensitive to the quality of the hash function or load factors, so it is often the choice when it is unknown how many and how frequently keys may be inserted or deleted from the hash table. On the other hand open addressing can be more space efficient as there are no linked lists. Linear probing has the advantage that it has small constants and works well with caches. But it suffers from primary clustering, which means its performance is sensitive to collisions and to high load factors. Quadratic probing, on the other hand, avoids primary clustering, but still suffers from secondary clustering and requires rehashing as soon as load factor reaches 50%. Although double hashing reduces clustering, so high load factors are possible, finding suitable pairs of hash functions is somewhat more difficult and increases the cost of a probe.

Next lecture, we will see that open addressing can be implemented in parallel.