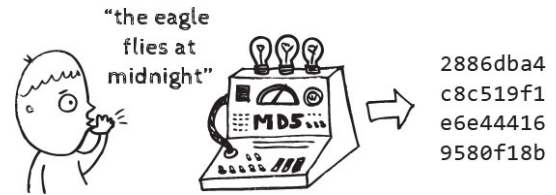


CSE 332

Data Structures and Parallelism

Hashing: Part 2

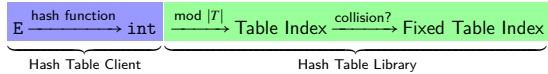


HashTable Review

1

Hash Tables

- Provides $\mathcal{O}(1)$ core Dictionary operations (**on average**)
- We call the key space the “universe”: U and the Hash Table T
- We should use this data structure **only** when we expect $|U| \gg |T|$
- (Or, the key space is non-integer values.)



Another Consideration?

What do we do when λ (the load factor) gets too large?

Hashing Choices

2

- Choose a hash function
- Choose a table size
- Choose a collision resolution strategy
 - Separate Chaining
 - Linear Probing
 - Quadratic Probing
 - Double Hashing
 - Other issues to consider:
- Choose an implementation of deletion
- Choose a λ that means the table is “too full”

We discussed the first few of these last time. We’ll discuss the rest today.

Review: Collisions

3

Definition (Collision)

A **collision** is when two distinct keys map to the same location in the hash table.

A good hash function attempts to avoid as many collisions as possible, but they are inevitable.

How do we deal with collisions?

There are multiple strategies:

- Separate Chaining
- Open Addressing
 - Linear Probing
 - Quadratic Probing
 - Double Hashing

Open Addressing

4

Definition (Open Addressing)

Open Addressing is a type of collision resolution strategy that resolves collisions by choosing a different location when the natural choice is full.

There are many types of open addressing. Here’s the key ideas:

- We **must** be able to duplicate the path we took.
- We want to use **all** the spaces in the table.
- We want to avoid putting lots of keys close together.

It turns out some of these are difficult to achieve. . .

Strategy #1: Linear Probing

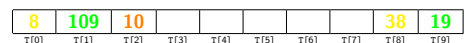
```

1 i = 0;
2 while (index in use) {
3   try (h(key) + i) % |T|
4 }

```

Example

Insert 38, 19, 8, 109, 10 into a hash table with hash function $h(x) = x$ and **linear probing**



(Items with the same hash code are the same color)

Open Addressing: Linear Probing

5

Strategy #1: Linear Probing

```

1 i = 0;
2 while (index in use) {
3   try (h(key) + i) % |T|
4 }
    
```

Example

Insert 38, 19, 8, 109, 10 into a hash table with hash function $h(x) = x$ and **linear probing**

8	109	10						38	19
$\tau(0)$	$\tau(1)$	$\tau(2)$	$\tau(3)$	$\tau(4)$	$\tau(5)$	$\tau(6)$	$\tau(7)$	$\tau(8)$	$\tau(9)$

(Items with the same hash code are the same color)

Other Operations with Linear Probing

- insert? Finds the **next** open spot. The worst case is $\mathcal{O}(n)$
- find? We have to retrace our steps. If the insert chain was k long, then $\text{find} \in \mathcal{O}(k)$.
- delete? We don't have a choice; we **must** use lazy deletion. What happens if we delete 19 and then do $\text{find}(109)$ in our example?

8	109	10						38	19
$\tau(0)$	$\tau(1)$	$\tau(2)$	$\tau(3)$	$\tau(4)$	$\tau(5)$	$\tau(6)$	$\tau(7)$	$\tau(8)$	$\tau(9)$

Analyzing Linear Probing

6

Which Criteria Does Linear Probing Meet?

- We want to use all the spaces in the table.
Yes! Linear probing will fill the whole table.
- We want to avoid putting lots of keys close together.
Uh... not so much

Primary Clustering

Primary Clustering is when different keys collide to form one big group.

8	109	10	101	20				38	19
$\tau(0)$	$\tau(1)$	$\tau(2)$	$\tau(3)$	$\tau(4)$	$\tau(5)$	$\tau(6)$	$\tau(7)$	$\tau(8)$	$\tau(9)$

Think of this as "clusters of many colors". Even though these keys are all different, they end up in a giant cluster.

In linear probing, we expect to get $\mathcal{O}(\lg n)$ size clusters.

This is really bad! But, how bad, really?

Analyzing Linear Probing

7

Load Factor & Space Usage

Note that $\lambda \leq 1$, and we will eventually get to $\lambda = 1$.

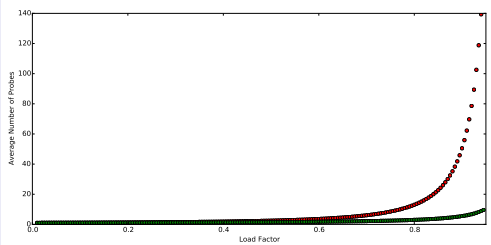
Average Number of Probes

Unsuccessful Search

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

Successful Search

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



Quadratic Probing

8

There's nothing theoretically wrong with open addressing that forces primary clustering. We'd like a different (easy to compute) function to probe with. That is:

Open Addressing In General

Choose a new function $f(x)$ and then probe with

$$(h(\text{key}) + f(i)) \bmod |T|$$

Strategy #2: Quadratic Probing

```

1 i = 0;
2 while (index in use) {
3   try (h(key) + i^2) % |T|
4 }
    
```

Example

Insert 89, 18, 49, 58, 79 into a hash table with hash function $h(x) = x$ and **quadratic probing**

49		58	79					18	89
$\tau(0)$	$\tau(1)$	$\tau(2)$	$\tau(3)$	$\tau(4)$	$\tau(5)$	$\tau(6)$	$\tau(7)$	$\tau(8)$	$\tau(9)$

$$\begin{aligned}
 h(58) \stackrel{i=0}{\rightarrow} 58 + 0^2 &= 0 \\
 \stackrel{i=1}{\rightarrow} 58 + 1^2 &= 0 \\
 \stackrel{i=2}{\rightarrow} 58 + 2^2 &= 2
 \end{aligned}$$

Another Quadratic Probing Example

9

Strategy #2: Quadratic Probing

```

1 i = 0;
2 while (index in use) {
3   try (h(key) + i^2) % |T|
4 }
    
```

Example

Insert 76, 40, 48, 5, 55, 47 into a hash table with hash function $h(x) = x$ and **quadratic probing**

48		5	55		40	76
$\tau(0)$	$\tau(1)$	$\tau(2)$	$\tau(3)$	$\tau(4)$	$\tau(5)$	$\tau(6)$

$$\begin{aligned}
 h(47) \stackrel{i=0}{\rightarrow} 47 + 0^2 &= 0 \\
 \stackrel{i=1}{\rightarrow} 47 + 1^2 &= 0 \\
 \stackrel{i=2}{\rightarrow} 47 + 2^2 &= 3 \\
 \stackrel{i=3}{\rightarrow} 47 + 3^2 &= 0 \\
 \stackrel{i=4}{\rightarrow} 47 + 4^2 &= 0 \\
 \stackrel{i=5}{\rightarrow} 47 + 5^2 &= 2
 \end{aligned}$$

We will never get a 1 or a 4!

This means we will never be able to insert 47. What's going on?

Quadratic Probing: Table Coverage

10

48		5	55		40	76
$\tau(0)$	$\tau(1)$	$\tau(2)$	$\tau(3)$	$\tau(4)$	$\tau(5)$	$\tau(6)$

Why Does insert(47) Fail?

For all i , $(5 + i^2) \bmod 7 \in \{0, 2, 5, 6\}$. The proof is by induction. This actually generalizes:

$$\text{For all } c, k, (c + i^2) \bmod k = (c + (i-k)^2) \bmod k$$

So, quadratic probing doesn't always **fill the table**.

The Good News!

If $|T|$ is prime and $\lambda < \frac{1}{2}$, then quadratic probing will find an empty slot in at most $\frac{|T|}{2}$ probes. So, if we keep $\lambda < \frac{1}{2}$, we don't need to detect cycles. The proof will be posted on the website.

So, does quadratic probing completely fix **clustering**?

Quadratic Probing: Clustering

11

With linear probing, we saw **primary clustering** (keys hashing near each other). Quadratic Probing fixes this by "jumping". Unfortunately, we still get **secondary clustering**:

Secondary Clustering

Secondary Clustering is when different keys hash to the same place and follow the same probing sequence.

39			29					9	19
$\tau_{(0)}$	$\tau_{(1)}$	$\tau_{(2)}$	$\tau_{(3)}$	$\tau_{(4)}$	$\tau_{(5)}$	$\tau_{(6)}$	$\tau_{(7)}$	$\tau_{(8)}$	$\tau_{(9)}$

Think of this as long probing chains of the same color. The keys all start at the same place; so, the chain gets really long.

We can avoid secondary clustering by using a probe function that **depends on the key**.

Linear/Quadratic/Double Hashing

12

Strategy #3: Double Hashing

```
1 i = 0;
2 while (index in use) {
3   try (h(key) + i*g(key)) % |T|
4 }
```

We insist $g(x) \neq 0$.

Example

Insert 13, 28, 33, 147, 43 into a hash table with:

$$\blacksquare h(x) = x$$

$$\blacksquare g(x) = 1 + \left(\frac{x}{|T|}\right) \bmod (|T| - 1)$$

using **double hashing**

			13				33	28	147
$\tau_{(0)}$	$\tau_{(1)}$	$\tau_{(2)}$	$\tau_{(3)}$	$\tau_{(4)}$	$\tau_{(5)}$	$\tau_{(6)}$	$\tau_{(7)}$	$\tau_{(8)}$	$\tau_{(9)}$

$$\begin{aligned} h(43) &= 43 \bmod 10 = 3 \\ g(43) &= 1 + \left(\frac{43}{10}\right) \bmod 9 = 3 \\ h(43) + g(43) &= 3 + 3 = 6 \\ h(43) + 2g(43) &= 3 + 2(3) = 9 \\ h(43) + 3g(43) &= 3 + 3(3) = 9 \\ &\vdots \\ h(43) + 4g(43) &= 3 + 4(3) = 15 \bmod 10 = 5 \end{aligned}$$

We got stuck again!

Double Hashing Analysis

13

Filling the Table

Just like with Quadratic Probing, we sometimes hit an infinite loop with double hashing. We will not get an infinite loop in the case with primes p, q such that $2 < q < p$:

- $h(\text{key}) = \text{key} \bmod p$
- $g(\text{key}) = q - (\text{key} \bmod q)$

Uniform Hashing

For double hashing, we assume **uniform hashing** which means:

$$\Pr[g(\text{key}1) \bmod p = g(\text{key}2) \bmod p] = \frac{1}{p}$$

Average Number of Probes

Unsuccessful Search

$$\frac{1}{1-\lambda}$$

Successful Search

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

This is way better than linear probing.

Where We Are

14

Separate Chaining is Easy!

- find, delete proportional to load factor on average
- insert can be constant if just push on front of list

Open Addressing is Tricky!

- Clustering issues
- Doesn't always use the whole table
- Why Use it?
 - Less memory allocation
 - Easier data representation

Now, let's move on to resizing the table.

Reshaping

15

When λ is too big, create a bigger table and copy over the items

When To Resize

- With separate chaining, we decide when to resize (should be $\lambda \leq 1$)
- With open addressing, we need to keep $\lambda < \frac{1}{2}$

New Table Size?

- Like always, we want around "twice as big"
- ... but it should still be prime
- So, choose the next prime about twice as big

How To Resize

Go through table, do standard insert for each into new table:

- Iterate over old table: $\mathcal{O}(n)$
- n inserts / calls to the hash function: $n \times \mathcal{O}(1) = \mathcal{O}(n)$
- But this is amortized $\mathcal{O}(1)$

Hashing and Comparing

16

A hash function isn't enough! We have to **compare** items:

- With separate chaining, we have to loop through the list checking if the item is what we're looking for
- With open addressing, we need to know when to stop probing

We have two options for this: **equality testing** or **comparison testing**.

- In Project 2, you will use both types.
- In Java, each Object has an equals method and a hashCode method

```
1 class Object {
2   boolean equals(Object o) {...}
3   int hashCode() {...}
4   ...
5 }
```

For any class, it **must be the case that**:

- If `a.equals(b)`, then `a.hashCode() == b.hashCode()`
- If `a.compareTo(b) == 0`, then `a.hashCode() == b.hashCode()`
- If `a.compareTo(b) < 0`, then `b.compareTo(a) > 0`
- If `a.compareTo(b) == 0`, then `b.compareTo(a) == 0`
- If `a.compareTo(b) < 0` and `b.compareTo(c) < 0`, then `a.compareTo(c) < 0`

```

1 int result = 17; // start at a prime
2 foreach field f
3 int fieldHashCode =
4     boolean: (f ? 1: 0)
5     byte, char, short, int: (int) f
6     long: (int) (f ^ (f >>> 32))
7     float: Float.floatToIntBits(f)
8     double: Double.doubleToLongBits(f), then above
9     Object: object.hashCode()
10 result = 31 * result + fieldHashCode;
11 return result;

```

- Hash Tables are one of the most important data structures
 - Efficient find, insert, and delete
 - based on sorted order are not so efficient
 - Useful in many, many real-world applications
 - Popular topic for job interview questions
- Important to use a good hash function
 - Good distribution, uses enough of keys values
 - Not overly expensive to calculate (bit shifts good!)
- Important to keep hash table at a good size
 - Prime Size
 - λ depends on type of table
- What we skipped: perfect hashing, universal hash functions, hopscotch hashing, cuckoo hashing