# Hashing
### Chapter 5 in Weiss

CSE 373

Data Structures and Algorithms

Ruth Anderson

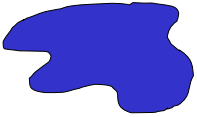11/03/2010                                                                 1

---

## Today's Outline

- **Announcements**
  - **Homework #4 due Fri, Nov 5 at the beginning of class**.

- **Today's Topics:**
  - **Disjoint Sets & Dynamic Equivalence**
  - **Hashing**

11/03/2010                                                                 2

---

## Hash Tables

hash table

- Constant time accesses!
- A **hash table** is an array of some fixed size, usually a prime number.
- General idea:

0

hash function:
**h(K)**

…

key space (e.g., integers, strings)          TableSize –1

11/03/2010                                                                 3

---

## Hash Tables

Key space of size M, but we only want to store subset of size N, where N<<M.

- Keys are identifiers in programs. Compiler keeps track of them in a symbol table.
- Keys are student names. We want to look up student records quickly by name.
- Keys are chess configurations in a chess playing program.
- Keys are URLs in a database of web pages.

11/03/2010                                                                 4

---

## Example

- key space = integers
- TableSize = 10

- $h(K) = K \bmod 10$

- **Insert**: 7, 18, 41, 94

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| 5 | |
| 6 | |
| 7 | |
| 8 | |
| 9 | |

11/03/2010                                                                 5

---

## Another Example

- key space = integers
- TableSize = 6

- $h(K) = K \bmod 6$

- **Insert**: 7, 18, 41, 34

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| 5 | |

11/03/2010                                                                 6

Student Activity

---

1

## Hash Functions

1. **simple/fast** to compute,
2. Avoid **collisions**
3. have keys distributed **evenly** among cells.

Perfect Hash function:

---

## Sample Hash Functions:

- key space = strings
- $s = s_0\, s_1\, s_2\, \ldots\, s_{k-1}$

1. $h(s) = s_0 \bmod \text{TableSize}$

2. $h(s) = \left( \displaystyle\sum_{i=0}^{k-1} s_i \right) \bmod \text{TableSize}$

3. $h(s) = \left( \displaystyle\sum_{i=0}^{k-1} s_i \cdot 37^i \right) \bmod \text{TableSize}$

---

## Designing a Hash Function for web URLs

$s = s_0\, s_1\, s_2\, \ldots\, s_{k-1}$

Issues to take into account:

$h(s) =$

**Student Activity**

---

## Collision Resolution

**Collision**: when two keys map to the same location in the hash table.

Two ways to resolve collisions:
1. Separate Chaining
2. Open Addressing (linear probing, quadratic probing, double hashing)

---

## Separate Chaining

**Insert**:
10
22
107
12
42

| 0 |  |
|---|---|
| 1 |  |
| 2 |  |
| 3 |  |
| 4 |  |
| 5 |  |
| 6 |  |
| 7 |  |
| 8 |  |
| 9 |  |

- **Separate chaining**: All keys that map to the same hash value are kept in a list ("bucket").

---

## Analysis of find

- The load factor, $\lambda$, of a hash table is the ratio:
$$\frac{N}{M} \quad \begin{array}{l} \leftarrow \text{no. of elements} \\ \leftarrow \text{table size} \end{array}$$

For separate chaining, $\lambda$ = average # of elements in a bucket

- unsuccessful:

- successful:

## How big should the hash table be?

- For Separate Chaining:

## tableSize: Why Prime?

- Suppose
  - data stored in hash table: 7160, 493, 60, 55, 321, 900, 810

  - tableSize = 10
    data hashes to 0, 3, 0, 5, 1, 0, 0

  - tableSize = 11
    data hashes to 10, 9, 5, 0, 2, 9, 7

> Real-life data tends to have a pattern
>
> Being a multiple of 11 is usually *not* the pattern ☺

## Open Addressing

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| 5 | |
| 6 | |
| 7 | |
| 8 | |
| 9 | |

**Insert**:
38
19
8
109
10

- **Linear Probing**: after checking spot h(k), try spot h(k)+1, if that is full, try h(k)+2, then h(k)+3, etc.

## Terminology Alert!

"**Open** Hashing"     "Closed Hashing"
equals     equals
"Separate Chaining"     "**Open** Addressing"

Weiss

## Linear Probing
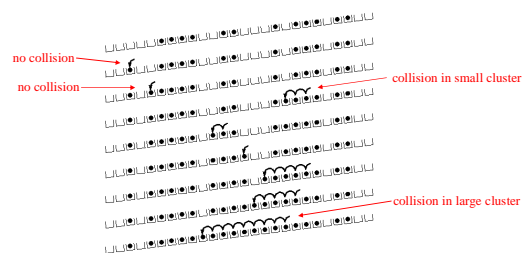
$$f(i) = i$$

- Probe sequence:
  - $0^{th}$ probe = h(k) mod TableSize
  - $1^{th}$ probe = (h(k) + 1) mod TableSize
  - $2^{th}$ probe = (h(k) + 2) mod TableSize
  - . . .
  - $i^{th}$ probe = (h(k) + i) mod TableSize

## Linear Probing – Clustering



no collision
no collision
collision in small cluster
collision in large cluster

[R. Sedgewick]

## Load Factor in Linear Probing

- For *any* $\lambda < 1$, linear probing *will* find an empty slot
- Expected # of probes (for large table sizes)
  - successful search:
  $$\frac{1}{2}\left(1 + \frac{1}{(1-\lambda)}\right)$$
  - unsuccessful search:
  $$\frac{1}{2}\left(1 + \frac{1}{(1-\lambda)^2}\right)$$
- Linear probing suffers from *primary clustering*
- Performance quickly degrades for $\lambda > 1/2$

11/03/2010

19

---

## Quadratic Probing

| Less likely to encounter Primary Clustering |
|---|

$$f(i) = i^2$$

- Probe sequence:
  - $0^{th}$ probe = h(k) mod TableSize
  - $1^{th}$ probe = (h(k) + 1) mod TableSize
  - $2^{th}$ probe = (h(k) + 4) mod TableSize
  - $3^{th}$ probe = (h(k) + 9) mod TableSize
  - . . .
  - $i^{th}$ probe = (h(k) + $i^2$) mod TableSize

11/03/2010

20

---

## Quadratic Probing

| 0 | |
|---|---|
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| 5 | |
| 6 | |
| 7 | |
| 8 | |
| 9 | |

Insert:
89
18
49
58
79

11/03/2010

21

---

## Quadratic Probing:

- h(k) = k mod 7
- Perform these inserts:
  - Insert(65)
  - Insert(10)
  - Insert(47)

| 0 | |
|---|---|
| 1 | |
| 2 | 93 |
| 3 | |
| 4 | |
| 5 | 40 |
| 6 | 76 |

11/03/2010

22

---

## Quadratic Probing Example

insert(76)   insert(40)   insert(48)   insert(5)   insert(55)
76%7 = 6     40%7 = 5     48%7 = 6    5%7 = 5     55%7 = 6

| 0 | |
|---|---|
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| 5 | |
| 6 | 76 |

But…  insert(47)
47%7 = 5

11/03/2010

23

---

## Quadratic Probing:
## Success guarantee for $\lambda < \frac{1}{2}$

- If size is prime and $\lambda < \frac{1}{2}$, then quadratic probing will find an empty slot in size/2 probes or fewer.
  - show for all `0 ≤ i,j ≤ size/2` and `i ≠ j`
    `(h(x) + i²) mod size ≠ (h(x) + j²) mod size`
  - by contradiction: suppose that for some `i ≠ j`:
    `(h(x) + i²) mod size = (h(x) + j²) mod size`
    $\Rightarrow$ `i² mod size = j² mod size`
    $\Rightarrow$ `(i² - j²) mod size = 0`
    $\Rightarrow$ `[(i + j)(i - j)] mod size = 0`
    BUT size does not divide `(i-j)` or `(i+j)`

11/03/2010

24

## Quadratic Probing: Properties

- For *any* λ < ½, quadratic probing will find an empty slot; for bigger λ, quadratic probing *may* find a slot

- Quadratic probing does not suffer from *primary* clustering: keys hashing to the same *area* are not bad

- But what about keys that hash to the same *spot*?
  - *Secondary Clustering!*

11/03/2010      25

---

## Double Hashing

$$f(i) = i * g(k)$$
where g is a second hash function

- Probe sequence:
  - $0^{th}$ probe = h(k) mod TableSize
  - $1^{th}$ probe = (h(k) + g(k)) mod TableSize
  - $2^{th}$ probe = (h(k) + 2*g(k)) mod TableSize
  - $3^{th}$ probe = (h(k) + 3*g(k)) mod TableSize
  - . . .
  - $i^{th}$ probe = (h($\underline{k}$) + i*g($\underline{k}$)) mod TableSize

11/03/2010      26

---

## Double Hashing Example

$i^{th}$ probe = (h($\underline{k}$) + i*g($\underline{k}$)) mod TableSize
h(k) = k mod 7 and g(k) = 5 – (k mod 5)

| | 76 | | 93 | | 40 | | 47 | | 10 | | 55 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | | 0 | | 0 | | 0 | | 0 | | 0 | |
| 1 | | 1 | | 1 | | 1 | 47 | 1 | 47 | 1 | 47 |
| 2 | | 2 | 93 | 2 | 93 | 2 | 93 | 2 | 93 | 2 | 93 |
| 3 | | 3 | | 3 | | 3 | | 3 | 10 | 3 | 10 |
| 4 | | 4 | | 4 | | 4 | | 4 | | 4 | 55 |
| 5 | | 5 | | 5 | 40 | 5 | 40 | 5 | 40 | 5 | 40 |
| 6 | 76 | 6 | 76 | 6 | 76 | 6 | 76 | 6 | 76 | 6 | 76 |

Probes   1      1      1      2      1      2

11/03/2010      27

---

## Resolving Collisions with Double Hashing

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| 5 | |
| 6 | |
| 7 | |
| 8 | |
| 9 | |

Hash Functions:
H(k) = k mod M
$H_2$(k) = 1 + ((k/M) mod (M-1))
M =

**Insert these values into the hash table in this order. Resolve any collisions with double hashing**:
13
28
33
147
43

11/03/2010      28

---

## Rehashing

**Idea**: When the table gets too full, create a bigger table (usually 2x as large) and hash all the items from the original table into the new table.

- When to rehash?
  - half full (λ = 0.5)
  - when an insertion fails
  - some other threshold
- Cost of rehashing?

11/03/2010      29

---

## Hashing Summary

- Hashing is one of the most important data structures.
- Hashing has many applications where operations are limited to find, insert, and delete.
- Dynamic hash tables have good amortized complexity.

11/03/2010      30