

Hash Tables

```
graph LR; 0((0)) --- E1[ ]; 1((1)) --> V1[025-612-0001]; 2((2)) --> V2[981-101-0002]; 3((3)) --- E2[ ]; 4((4)) --> V3[451-229-0004];
```

© 2013 Goodrich, Tamassia, Goldwasser

Hash Tables

1

Recall the notion of a Map

- Intuitively, a map M supports the abstraction of using keys as indices with a syntax such as $M[k]$.
- As a mental warm-up, consider a restricted setting in which a map with n items uses keys that are known to be integers in a range from 0 to $N - 1$, for some $N \geq n$.

0	1	2	3	4	5	6	7	8	9	10
	D		Z			C	Q			

© 2013 Goodrich, Tamassia, Goldwasser

Hash Tables

2

More General Kinds of Keys

- But what should we do if our keys are not integers in the range from 0 to $N - 1$?
 - Use a **hash function** to map general keys to corresponding indices in a table.
 - For instance, the last four digits of a Social Security number.

```
graph LR; 0((0)) --- E1[ ]; 1((1)) --> V1[025-612-0001]; 2((2)) --> V2[981-101-0002]; 3((3)) --- E2[ ]; 4((4)) --> V3[451-229-0004];
```

© 2013 Goodrich, Tamassia, Goldwasser

Hash Tables

3

Hash Functions and Hash Tables

- A **hash function** h maps keys of a given type to integers in a fixed interval $[0, N - 1]$
- Example:

$$h(x) = x \bmod N$$
is a hash function for integer keys
- The integer $h(x)$ is called the **hash value** of key x .
- A **hash table** for a given key type consists of
 - Hash function h
 - Array (called table) of size N
- When implementing a map with a hash table, the goal is to store item (k, o) at index $i = h(k)$

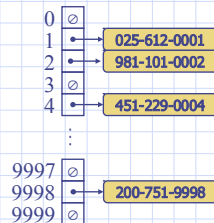
© 2013 Goodrich, Tamassia, Goldwasser

Hash Tables

4

SSN Example

- We design a hash table for a map storing entries as (SSN, Name), where SSN (social security number) is a nine-digit positive integer
- Our hash table uses an array of size $N = 10,000$ and the hash function $h(x) =$ last four digits of x



Hash Functions



- A hash function is usually specified as the composition of two functions:
 - Hash code:**
 $h_1: \text{keys} \rightarrow \text{integers}$
 - Compression function:**
 $h_2: \text{integers} \rightarrow [0, N - 1]$
- The hash code is applied first, and the compression function is applied next on the result, i.e.,
$$h(x) = h_2(h_1(x))$$
- The goal of the hash function is to “disperse” the keys in an apparently random way

Hash Codes



- **Memory address:**
 - We reinterpret the memory address of the key object as an integer Good in general, except for numeric and string keys
- **Integer cast:**
 - We reinterpret the bits of the key as an integer
 - Suitable for keys of length less than or equal to the number of bits of the integer
- **Component sum:**
 - We partition the bits of the key into components of fixed length (e.g., 16 or 32 bits) and we sum the components (ignoring overflows)
 - Suitable for numeric keys of fixed length greater than or equal to the number of bits of the integer type

Hash Codes (cont.)

- Polynomial accumulation:
 - We partition the bits of the key into a sequence of components of fixed length (e.g., 8, 16 or 32 bits)

$$a_0 a_1 \dots a_{n-1}$$
 - We evaluate the polynomial

$$p(z) = a_0 + a_1 z + a_2 z^2 + \dots + a_{n-1} z^{n-1}$$
 at a fixed value z , ignoring overflows
 - Especially suitable for strings (e.g., the choice $z = 33$ gives at most 6 collisions on a set of 50,000 English words)
- Polynomial $p(z)$ can be evaluated in $O(n)$ time using Horner's rule:
 - The following polynomials are successively computed, each from the previous one in $O(1)$ time

$$p_0(z) = a_{n-1}$$

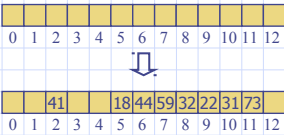
$$p_i(z) = a_{n-i-1} + z p_{i-1}(z)$$
 ($i = 1, 2, \dots, n-1$)
 - We have $p(z) = p_{n-1}(z)$

Hash Table with Separate Chaining

```
1 class ChainHashMap(HashMapBase):
2     """Hash map implemented with separate chaining for collision resolution."""
3
4     def _bucket_getitem(self, j, k):
5         bucket = self._table[j]
6         if bucket is None:
7             raise KeyError('Key Error: ' + repr(k)) # no match found
8         return bucket[k] # may raise KeyError
9
10    def _bucket_setitem(self, j, k, v):
11        if self._table[j] is None: # bucket is new to the table
12            self._table[j] = UnsortedTableMap()
13            oldsize = len(self._table[j])
14            self._table[j][k] = v # key was new to the table
15            if len(self._table[j]) > oldsize: # increase overall map size
16                self._n += 1
17
18    def _bucket_delitem(self, j, k):
19        bucket = self._table[j]
20        if bucket is None:
21            raise KeyError('Key Error: ' + repr(k)) # no match found
22        del bucket[k] # may raise KeyError
23
24    def __iter__(self):
25        for bucket in self._table:
26            if bucket is not None: # a nonempty slot
27                for key in bucket:
28                    yield key
```

Linear Probing

- Open addressing: the colliding item is placed in a different cell of the table
 - Linear probing: handles collisions by placing the colliding item in the next (circularly) available table cell
 - Each table cell inspected is referred to as a "probe"
 - Colliding items lump together, causing future collisions to cause a longer sequence of probes
- Example:
- $h(x) = x \bmod 13$
 - Insert keys 18, 41, 22, 44, 59, 32, 31, 73, in this order



Search with Linear Probing



- Consider a hash table A that uses linear probing
- $get(k)$
 - We start at cell $h(k)$
 - We probe consecutive locations until one of the following occurs
 - An item with key k is found, or
 - An empty cell is found, or
 - N cells have been unsuccessfully probed

```
Algorithm get(k)
    i ← h(k)
    p ← 0
    repeat
        c ← A[i]
        if c = ∅
            return null
        else if c.getKey() = k
            return c.getValue()
        else
            i ← (i + 1) mod N
            p ← p + 1
    until p = N
    return null
```

Updates with Linear Probing

- To handle insertions and deletions, we introduce a special object, called *AVAILABLE*, which replaces deleted elements
- $put(k, o)$
 - We throw an exception if the table is full
 - We start at cell $h(k)$
 - We probe consecutive cells until one of the following occurs
 - A cell i is found that is either empty or stores *AVAILABLE*, or
 - N cells have been unsuccessfully probed
 - We store (k, o) in cell i
- $remove(k)$
 - We search for an entry with key k
 - If such an entry (k, o) is found, we replace it with the special item *AVAILABLE* and we return element o
 - Else, we return *null*

Hash Table with Linear Probing

```
class ProbHashMap(HashMapBase):
    """Hash map implemented with linear probing for collision resolution"""
    _AVAIL = object() # sentinel marks locations of previous deletions

    def __init__(self, j):
        """Return True if index j is available in table"""
        return self._table[j] is None or self._table[j] is ProbHashMap._AVAIL

    def _find_slot(self, j, k):
        """Search for key k in bucket at index j"""
        # Return (success, index) tuple, described as follows:
        # If match was found, success is True and index denotes its location.
        # If no match found, success is False and index denotes first available slot.
        ...

    firstAvail = None
    while True:
        if self._is_available(j):
            firstAvail = j
            if firstAvail is None:
                return (False, firstAvail)
            elif k == self._table[j].key:
                return (True, j)
            j = (j + 1) % len(self._table) # keep looking (cyclically)

    def _bucket_getitem(self, j, k):
        found, s = self._find_slot(j, k)
        if not found:
            raise KeyError('Key Error: ' + repr(k)) # no match found
        return self._table[s].value

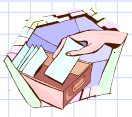
    def _bucket_setitem(self, j, k, v):
        found, s = self._find_slot(j, k)
        if not found:
            self._table[s] = self._Item(k, v) # insert new item
            self._n += 1 # size has increased
        else:
            self._table[s].value = v # overwrite existing

    def _bucket_delitem(self, j, k):
        found, s = self._find_slot(j, k)
        if not found:
            raise KeyError('Key Error: ' + repr(k)) # no match found
        self._table[s] = ProbHashMap._AVAIL # mark as vacated

    def __iter__(self):
        for j in range(len(self._table)):
            if not self._is_available(j):
                yield self._table[j].key
```

Double Hashing

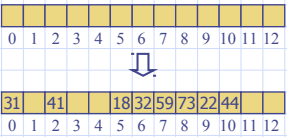
- Double hashing uses a secondary hash function $d(k)$ and handles collisions by placing an item in the first available cell of the series $(i + jd(k)) \bmod N$ for $j = 0, 1, \dots, N - 1$
- The secondary hash function $d(k)$ cannot have zero values
- The table size N must be a prime to allow probing of all the cells
- Common choice of compression function for the secondary hash function:
$$d_2(k) = q - k \bmod q$$
where
 - $q < N$
 - q is a prime
- The possible values for $d_2(k)$ are $1, 2, \dots, q$



Example of Double Hashing

- Consider a hash table storing integer keys that handles collision with double hashing
 - $N = 13$
 - $h(k) = k \bmod 13$
 - $d(k) = 7 - k \bmod 7$
- Insert keys 18, 41, 22, 44, 59, 32, 31, 73, in this order

k	h(k)	d(k)	Probes
18	5	3	5
41	2	1	2
22	9	6	9
44	5	5	5 10
59	7	4	7
32	6	3	6
31	5	4	5 9 0
73	8	4	8



Performance of Hashing

- In the worst case, searches, insertions and removals on a hash table take $O(n)$ time
- The worst case occurs when all the keys inserted into the map collide
- The load factor $\alpha = n/N$ affects the performance of a hash table
- Assuming that the hash values are like random numbers, it can be shown that the expected number of probes for an insertion with open addressing is $1 / (1 - \alpha)$
- The expected running time of all the dictionary ADT operations in a hash table is $O(1)$
- In practice, hashing is very fast provided the load factor is not close to 100%
- Applications of hash tables:
 - small databases
 - compilers
 - browser caches

