

Hashing

Hash Tables


- *Hash table* ADT which supports only a subset of the operations allowed by binary search trees.
- The implementation of hash tables is called **hashing**.
- Hashing is a technique used for performing insertions, deletions and finds in constant average time (i.e. $O(1)$)
- This data structure, however, is not efficient in operations that require any ordering information among the elements, such as findMin, findMax and printing the entire table in sorted order.

General Idea

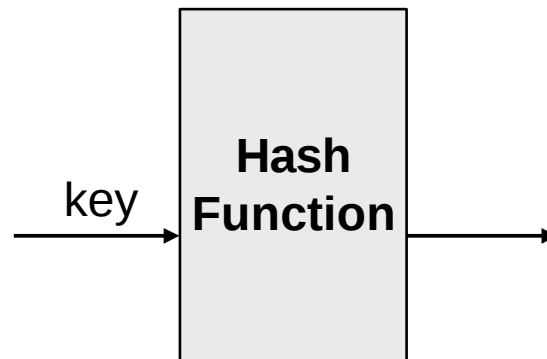
- The ideal hash table structure is merely an array of some fixed size, containing the items.
- A stored item needs to have a data member, called **key**, that will be used in computing the index value for the item.
 - Key could be an *integer*, a *string*, etc
 - e.g. a name or Id that is a part of a large employee structure
- The size of the array is *TableSize*.
- The items that are stored in the hash table are indexed by values from 0 to $TableSize - 1$.
- Each key is mapped into some number in the range 0 to $TableSize - 1$.
- The mapping is called a *hash function*.
- Number of keys or elements is N

Example

Items
john 25000
Phil 31250
dave 27500
mary 28200



key



Hash	
0	
1	
2	
3	
4	john 25000
5	phil 31250
6	
7	
8	dave 27500
9	mary 28200

Hash Function

- The hash function (map keys to indices):
 - must be fast and easy to compute.
 - equal keys should yield equal hash codes
 - must distribute the keys evenly among the cells.
 - Be in range 0 to $TableSize - 1$.
- Hash Function consists of :
 - Hash Code map (keys to integer values)
 - +
 - Compression map (next to indices of table)

Hash function

Problems:

- Keys may not be numeric.
- Number of possible keys is much larger than the space available in table.
- Different keys may map into same location
 - Hash function is not one-to-one => collision.
 - If there are too many collisions, the performance of the hash table will suffer dramatically.

Hash Functions

- If the input keys are integers then simply $Key \bmod TableSize$ is a general strategy.
 - Unless key happens to have some undesirable properties. (e.g. all keys end in 0 and we use mod 10)
- If the keys are strings, hash function needs more care.
 - First convert it into a numeric value (hash code).

Binning

- Key Range 0 to 999, uniformly distributed
- TableSize : 10
- `int hash(int i) {return i % 100;}`
 - 0-99 will hash to 0
 - 101-199 will hash to 1 and so on
- In general for any TableSize :
 - $\text{Key}/(1000/\text{TableSize})$.

Simple Hash Function

- Now assume TableSize = 100 and keys are distributed in wide range.

```
int hash(int i) {return i % 100;}
```

keys
00005, 00010, ... 00100
00105, 00110, ... 00200
00205, 00210, ... 00300
00305, 00310, ... 00400
⋮

Power of prime

- preserves the uniform distribution if already present in a set of keys
- tends to increase the uniformity in a set of keys if TableSize is a prime number
- So try to use prime number for the TableSize
- `int hash (int i) { return i % 101;}`

keys	hash to
00005, 00010, ... 00100	5, 10, ... , 100
00105, 00110, ... 00200	4, 9, ... , 99
00205, 00210, ... 00300	3, 8, ... , 98
00305, 00310, ... 00400	2, 7, ... , 97
⋮	⋮

Some Hash Code methods

- **Truncation:**
 - e.g. 123456789 map to a table of 1000 addresses by picking 3 digits of the key.
- **Folding:**
 - e.g. 123|456|789: add them and take mod.
- **Key mod TableSize:**
 - size of the table, better if it is prime.
- **Squaring and mid-Squaring:**
 - Square the key and then truncate
- **Radix conversion:**
 - e.g. 1 2 3 4 treat it to be base 11, truncate if necessary.

Folding

Shift Folding

Key (as integer) is divided into segments, each segment except possibly the last having the same number of digits. These segments are then added to obtain the index (or mod TableSize)

Key = 76123451001214 with segments of size 3 digits. The segments for our key are 761, 234, 510, 012, and 14.

The index is $761 + 234 + 510 + 012 + 14 = 1531$.

Folding at Boundaries :, the digits in alternate segments are reversed before adding

Segments after digit reversal are 761, 432, 510, 210, and 14;
The index is $761 + 432 + 510 + 210 + 14 = 1927$.

Hash Function 1

- Add up the ASCII values of all characters of the key.

```
int hash(const string &key, int tableSize)
{
    int hasVal = 0;

    for (int i = 0; i < key.length(); i++)
        hashVal += key[i];
    return hashVal % tableSize;
}
```

- Simple to implement and fast.
- Words differing by transposing of chars have same value.
- If table size is large, the does not distribute the keys well.
 - e.g. Table size =10000, key length <= 8, the hash function can assume values only between 0 and 1016

Hash Function 2

- Examine only the first 3 characters of the key.

```
int hash (const string &key, int tableSize)
{
    return (key[0]+27 * key[1] + 729*key[2]) %
    tableSize;
}
```

- In theory, $26 * 26 * 26 = 17576$ different words can be generated. However, English is not random, only **2851** different combinations are possible.
- Thus, this function although easily computable, is also not appropriate if the hash table is reasonably large.

Hash Function 3

$$\text{hash}(\text{key}) = \sum_{i=0}^{\text{KeySize}-1} \text{Key}[\text{KeySize}-i-1] \cdot 37^i$$

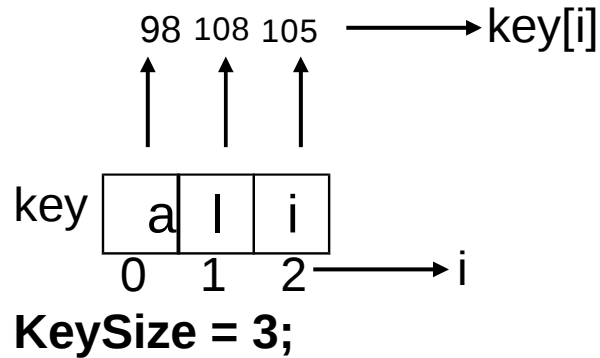
```
int hash (const string &key, int tableSize)
{
    int hashVal = 0;

    for (int i = 0; i < key.length(); i++) hashVal
        = 37 * hashVal + key[i];

    hashVal %=tableSize;
    if (hashVal < 0) /* in case overflows occurs */
        hashVal += tableSize;

    return hashVal;
};
```

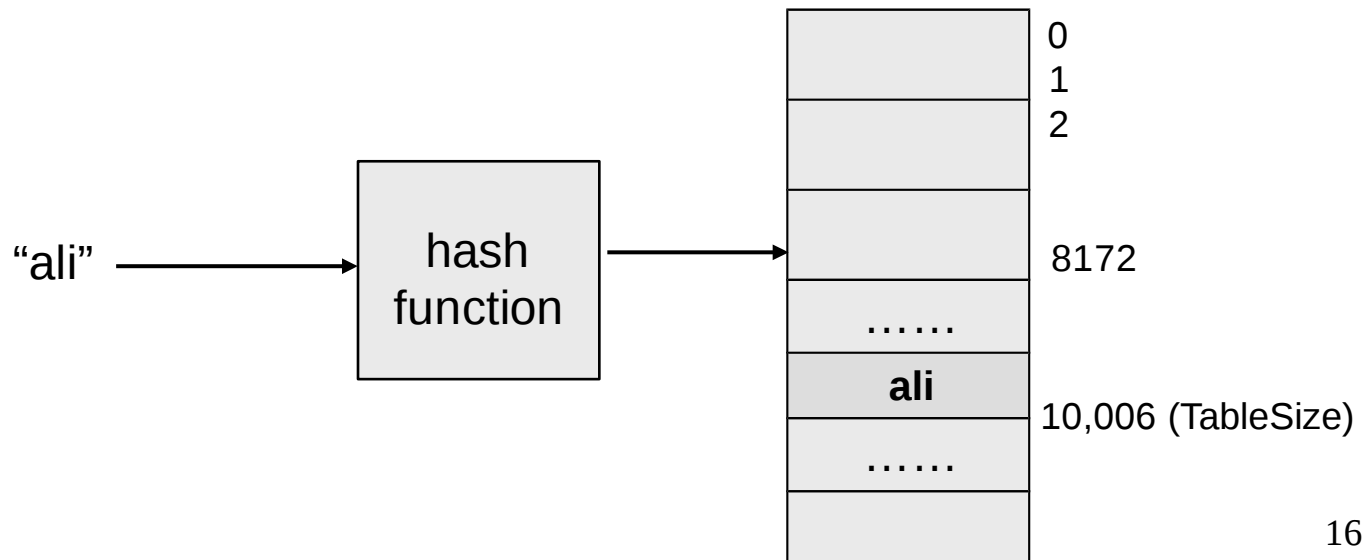
Hash function for strings:



Step 1: Reverse the given string

Step 2: Compute hashvalue

$$\text{hash}(\text{"ali"}) = (105 * 1 + 108 * 37 + 98 * 37^2) \% 10,007 = 8172$$



Compression Functions

Compress *hash(key)* or Hash Code to bucket size N

- **The Division Method**

$$|i| \bmod N, (N \text{ prime})$$

Eg. $N=100$ with Hash Codes (200, 205, 210 600)

- **The MOD Method**

$$|a_i + b| \bmod N \quad (a \bmod N \text{ not } 0)$$

Collision Resolution

- If, when an element is inserted, it hashes to the same value as an already inserted element, then we have a collision and need to resolve it.
- There are several methods for dealing with this:
 - **Separate chaining**
 - **Open addressing**
 - Linear Probing
 - Quadratic Probing
 - Double Hashing

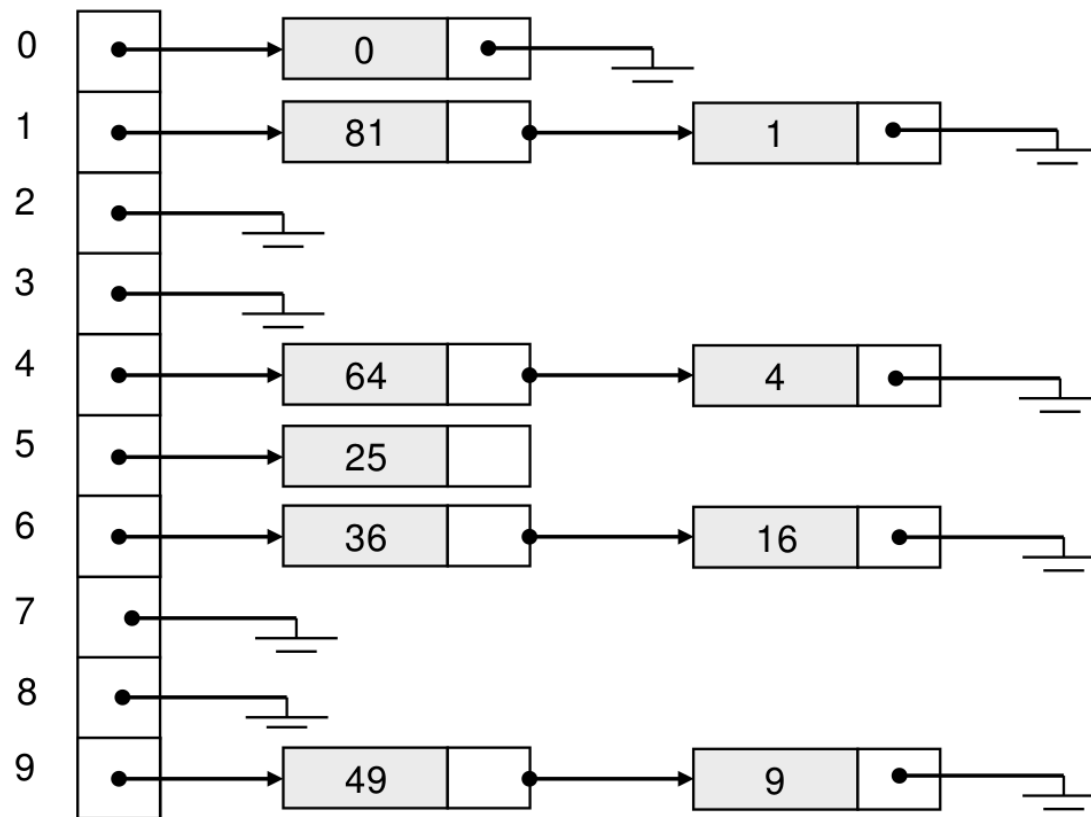
Separate Chaining

- The idea is to keep a list of all elements that hash to the same value.
 - The array elements are pointers to the first nodes of the lists.
 - A new item is inserted to the front of the list.
- Advantages:
 - Better space utilization for large items.
 - Simple collision handling: searching linked list.
 - Overflow: we can store more items than the hash table size.
 - Deletion is quick and easy: deletion from the linked list.

Example

Keys: 0, 1, 4, 9, 16, 25, 36, 49, 64, 81

$\text{hash}(\text{key}) = \text{key} \% 10.$



Operations

- **Initialization:** all entries are set to NULL
- **Find:**
 - locate the cell using hash function.
 - sequential search on the linked list in that cell.
- **Insertion:**
 - Locate the cell using hash function.
 - (If the item does not exist) insert it as the first item in the list.
- **Deletion:**
 - Locate the cell using hash function.
 - Delete the item from the linked list.

Hash Table Class for separate chaining

```
template <class HashedObj>  class HashTable
{
    public:
        HashTable(const HashedObj & notFound, int size=101 );
        HashTable( const HashTable & rhs )
            :ITEM_NOT_FOUND( rhs.ITEM_NOT_FOUND ),
              theLists( rhs.theLists ) { }

        const HashedObj & find( const HashedObj & x ) const;

        void makeEmpty( );
        void insert( const HashedObj & x );  void remove( const
        HashedObj & x );

        const HashTable & operator=( const HashTable & rhs );
    private:
        vector<List<HashedObj> > theLists; // The array of
        Lists
        const HashedObj ITEM_NOT_FOUND;
};

int hash( const string & key, int tableSize );  int
hash( int key, int tableSize );
```

Insert routine

```
/**
 * Insert item x into the hash table. If the item is
 * already present, then do nothing.
 */
template <class HashedObj>
void HashTable<HashedObj>::insert(const HashedObj &
x )
{
    List<HashedObj> & whichList = theLists[ hash( x,
                                                    theLists.size( ) ) ];
    ListItr<HashedObj> itr = whichList.find( x );
                                if( !itr.isValid( ) )
                                    whichList.insert( x,
                                                    whichList.zeroth( ) );
}
```

Remove routine

```
/**
 * Remove item x from the hash table.
 */
template <class HashedObj>
void HashTable<HashedObj>::remove( const HashedObj &
x )
{
    theLists[hash(x, theLists.size())].remove( x );
}
```


Find routine

```
/**
 * Find item x in the hash table.
 * Return the matching item or ITEM_NOT_FOUND if not found
 */
template <class HashedObj>
const HashedObj & HashTable<HashedObj>::find( const
    HashedObj & x ) const
{
    ListItr<HashedObj> itr;
    itr = theLists[ hash( x, theLists.size( ) ) ].find( x );
    if(!itr.isValid())
        return ITEM_NOT_FOUND;  else
        return itr.retrieve( );
}
```

Analysis of Separate Chaining

- Collisions are very likely.
 - How likely and what is the average length of lists?
- Load factor λ definition:
 - Ratio of number of elements (N) in a hash table to the hash *TableSize*.
 - i.e. $\lambda = N/TableSize$
 - The average length of a list is also λ .
 - For chaining λ is not bound by 1; it can be > 1 .

Cost of searching

- **Cost** = Constant time to evaluate the hash function + time to traverse the list.
- **Unsuccessful search:**
 - We have to traverse the entire list, so we need to compare λ nodes on the average.
- **Successful search:**
 - List contains the one node that stores the searched item + 0 or more other nodes.
 - Expected # of other nodes = $x = (N-1)/M$ which is essentially λ , since M is presumed large.
 - On the average, we need to check *half* of the *other nodes* while searching for a certain element
 - Thus average search cost = $1 + \lambda/2$

Summary

- The analysis shows us that the table size is not really important, but the load factor is.
- TableSize should be as *large* as the number of expected elements in the hash table.
 - To keep load factor around 1.
- TableSize should be *prime* for even distribution of keys to hash table cells.

Hashing: Open Addressing

Collision Resolution with Open Addressing

- Separate chaining has the disadvantage of using linked lists.
 - Requires the implementation of a second data structure.
- In an open addressing hashing system, all the data go inside the table.
 - Thus, a bigger table is needed.
 - Generally the load factor should be below 0.5.
 - If a collision occurs, alternative cells are tried until an empty cell is found.

Open Addressing

- More formally:
 - Cells $h_0(x)$, $h_1(x)$, $h_2(x)$, ... are tried in succession where
$$h_i(x) = (\text{hash}(x) + f(i)) \bmod \text{TableSize}, \text{ with } f(0) = 0.$$
 - The function f is the collision resolution strategy.
- There are three common collision resolution strategies:
 - Linear Probing $\{f(i) = i\}$
 - Quadratic probing $\{f(i) = i^2\}$
 - Double hashing $\{f(i) = i * \text{hash}_2(x)\}$

Linear Probing

- In linear probing, collisions are resolved by sequentially scanning an array (with wraparound) until an empty cell is found.
 - i.e. f is a linear function of i , typically $f(i) = i$.
- Example:
 - Insert items with keys: 89, 18, 49, 58, 9 into an empty hash table.
 - Table size is 10.
 - Hash function is $\text{hash}(x) = x \bmod 10$.
 - $f(i) = i$;

Figure 20.4

Linear probing
hash table after
each insertion

hash (89, 10) = 9
hash (18, 10) = 8
hash (49, 10) = 9
hash (58, 10) = 8
hash (9, 10) = 9

	<i>After insert 89</i>	<i>After insert 18</i>	<i>After insert 49</i>	<i>After insert 58</i>	<i>After insert 9</i>
0			49	49	49
1				58	58
2					9
3					
4					
5					
6					
7					
8		18	18	18	18
9	89	89	89	89	89

Figure 20.6

A quadratic
probing hash table
after each
insertion (note that
the table size was
poorly chosen
because it is not a
prime number).

hash (89, 10) = 9
hash (18, 10) = 8
hash (49, 10) = 9
hash (58, 10) = 8
hash (9, 10) = 9

	<i>After insert 89</i>	<i>After insert 18</i>	<i>After insert 49</i>	<i>After insert 58</i>	<i>After insert 9</i>
0			49	49	49
1					
2				58	58
3					9
4					
5					
6					
7					
8		18	18	18	18
9	89	89	89	89	89

Linear Probing

Insert: A_5, A_2, A_3

0	
1	
2	A_2
3	A_3
4	
5	A_5
6	
7	
8	
9	

(a)

B_5, A_9, B_2

0	
1	
2	A_2
3	A_3
4	B_2
5	A_5
6	B_5
7	
8	
9	A_9

(b)

B_9, C_2

0	B_9
1	
2	A_2
3	A_3
4	B_2
5	A_5
6	B_5
7	C_2
8	
9	A_9

(c)

Clustering Problem

- As long as table is big enough, a free cell can always be found, but the time to do so can get quite large.
- Worse, even if the table is relatively empty, blocks of occupied cells start forming.
- This effect is known as *primary clustering*.
- Any key that hashes into the cluster will require several attempts to resolve the collision, and then it will add to the cluster.

Quadratic Probing

- Quadratic Probing eliminates primary clustering problem of linear probing.
- Collision function is quadratic.
 - The popular choice is $f(i) = i^2$.
- If the hash function evaluates to h and a search in cell h is inconclusive, we try cells $h + 1^2, h + 2^2, \dots, h + i^2$.
 - i.e. It examines cells 1,4,9 and so on away from the original probe.
- Remember that subsequent probe points are a quadratic number of positions from the *original probe point*.

Figure 20.6

A quadratic
probing hash table
after each
insertion (note that
the table size was
poorly chosen
because it is not a
prime number).

hash (89, 10) = 9
hash (18, 10) = 8
hash (49, 10) = 9
hash (58, 10) = 8
hash (9, 10) = 9

	<i>After insert 89</i>	<i>After insert 18</i>	<i>After insert 49</i>	<i>After insert 58</i>	<i>After insert 9</i>
0			49	49	49
1					
2				58	58
3					9
4					
5					
6					
7					
8		18	18	18	18
9	89	89	89	89	89

Quadratic Probing

$$h(K) + i^2, h(K) - i^2 \text{ for } i = 1, 2, \dots, (TSize - 1)/2$$

Including the first attempt to hash K , this results in the sequence:

$$h(K), h(K) + 1, h(K) - 1, h(K) + 4, h(K) - 4, \dots, h(K) + (TSize - 1)^2/4, \\ h(K) - (TSize - 1)^2/4$$

all divided modulo $TSize$. The size of the table should not be an even number, because only the even positions or only the odd positions are tried depending on the value of $h(K)$. Ideally, the table size should be a prime $4j + 3$ of an integer j , which guarantees the inclusion of all positions in the probing sequence (Radke, 1970). For example, if $j = 4$, then $TSize = 19$, and assuming that $h(K) = 9$ for some K , the resulting sequence of probes is¹

$$9, 10, 8, 13, 5, 18, 0, 6, 12, 15, 3, 7, 11, 1, 17, 16, 2, 14, 4$$

Quadratic Probing

- Problem:
 - We may not be sure that we will probe all locations in the table (i.e. there is no guarantee to find an empty cell if table is more than half full.)
 - If the hash table size is not prime this problem will be much severe.
- However, there is a theorem stating that:
 - If the table size is *prime* and load factor is not larger than 0.5, all probes will be to different locations and an item can always be inserted.

Find and Delete

- The find algorithm follows the same probe sequence as the insert algorithm.
 - A find for 58 would involve 4 probes.
 - A find for 19 would involve 5 probes.
- We must use *lazy deletion* (i.e. marking items as deleted)
 - Standard deletion (i.e. physically removing the item) cannot be performed.
 - e.g. remove 89 from hash table.

Figure 20.6

A quadratic
probing hash table
after each
insertion (note that
the table size was
poorly chosen
because it is not a
prime number).

hash (89, 10) = 9
hash (18, 10) = 8
hash (49, 10) = 9
hash (58, 10) = 8
hash (9, 10) = 9

	<i>After insert 89</i>	<i>After insert 18</i>	<i>After insert 49</i>	<i>After insert 58</i>	<i>After insert 9</i>
0			49	49	49
1					
2				58	58
3					9
4					
5					
6					
7					
8		18	18	18	18
9	89	89	89	89	89

Theorem

- If quadratic probing is used, and the table size is prime, then a new element can always be inserted if the table is at least half empty.

Proof

- Let M be the size of the table and it is *prime*. We show that the first $\lceil M/2 \rceil$ alternative locations are distinct.
- Let two of these locations are $h + i^2$ and $h + j^2$, where i, j are two probes s.t. $0 \leq i, j \leq \lfloor M/2 \rfloor$. Suppose for the sake of contradiction, that these two locations are the same but $i \neq j$. Then

$$\begin{aligned} h + i^2 &= h + j^2 \pmod{M} & i^2 &= j^2 \pmod{M} \\ i^2 - j^2 &= 0 \pmod{M} & (i-j)(i+j) &= 0 \pmod{M} \end{aligned}$$

- Because M is prime, either $(i-j)$ or $(i+j)$ is divisible by M . Neither of these possibilities can occur. Thus we obtain a contradiction.
- It follows that the first $\lceil M/2 \rceil$ alternative are all distinct and since there are at most $\lfloor M/2 \rfloor$ items in the hash table it is guaranteed that an insertion must succeed if the table is at least half full.

Some considerations

- How efficient is calculating the quadratic probes?
 - Linear probing is easily implemented.
Quadratic probing appears to require * and % operations.
 - However by the use of the following trick, this is overcome:
 - $H_i = H_{i-1} + 2i - 1 \pmod{M}$

Some Considerations

- What happens if load factor gets too high?
 - Dynamically expand the table as soon as the load factor reaches 0.5, which is called *rehashing*.
 - Always double to a prime number.
 - When expanding the hash table, reinsert the new table by using the new hash function.

Analysis of Quadratic Probing

- Quadratic probing has not yet been mathematically analyzed.
- Although quadratic probing eliminates primary clustering, elements that hash to the same location will probe the same alternative cells. This is known as *secondary clustering*.
- Techniques that eliminate secondary clustering are available.
 - the most popular is *double hashing*.

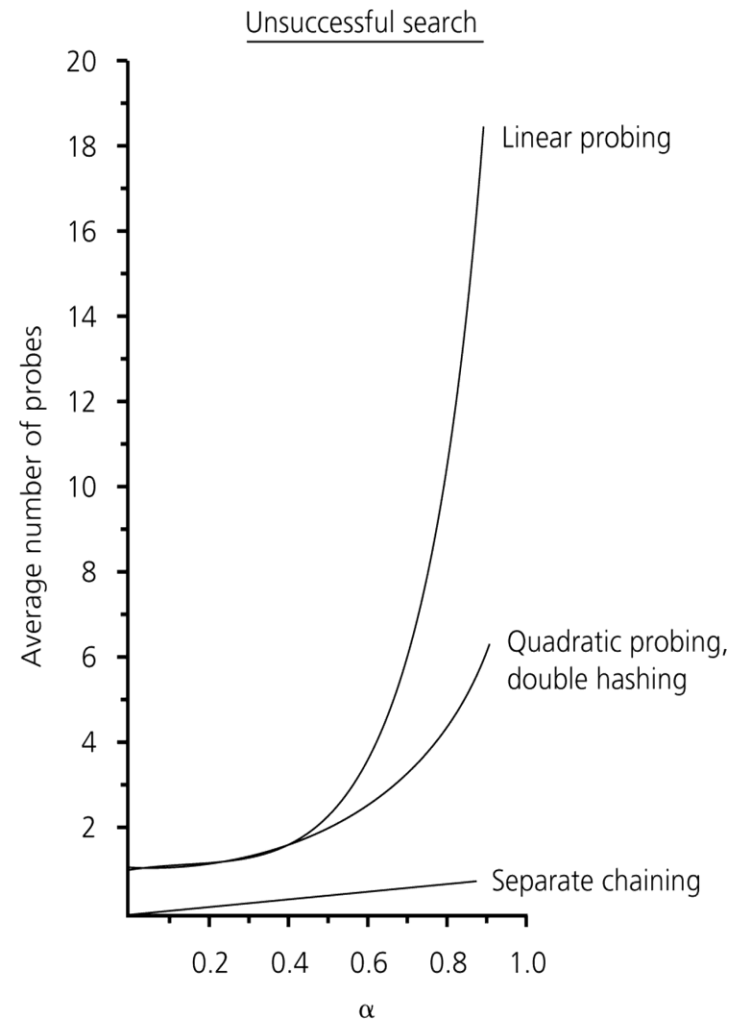
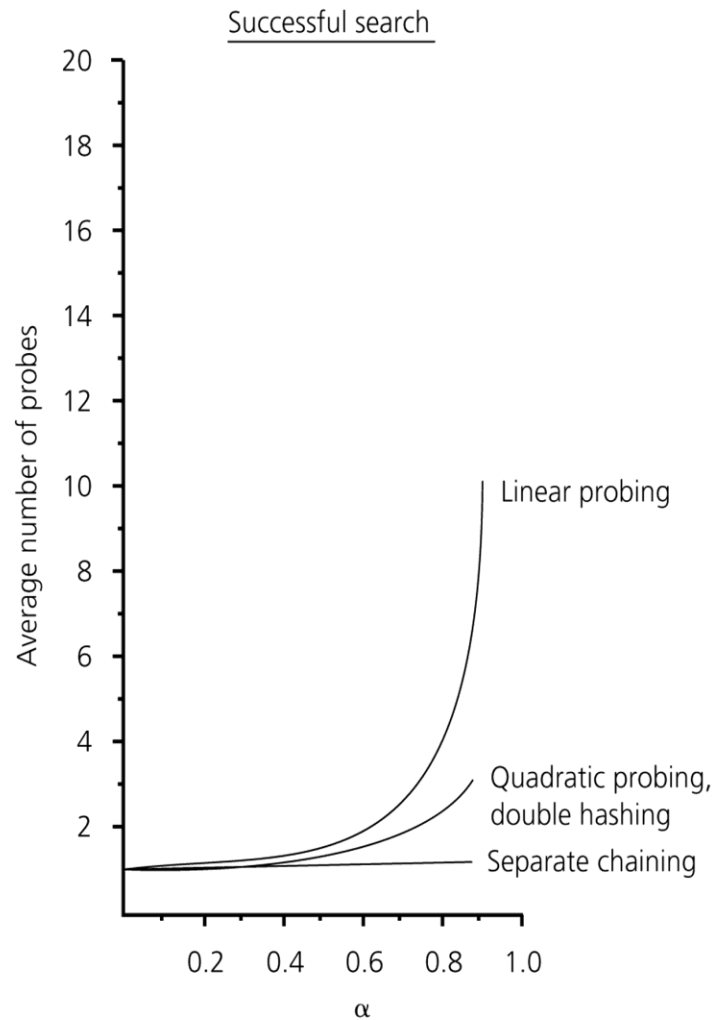
Double Hashing

- A second hash function is used to drive the collision resolution.
 - $f(i) = i * hash_2(x)$
- We apply a second hash function to x and probe at a distance $hash_2(x)$, $2*hash_2(x)$, ... and so on.
- The function $hash_2(x)$ must never evaluate to zero.
 - e.g. Let $hash_2(x) = x \bmod 9$ and try to insert 99 in the previous example.
- A function such as $hash_2(x) = R - (x \bmod R)$ with R a prime smaller than TableSize will work well.
 - e.g. try $R = 7$ for the previous example. $(7 - x \bmod 7)$

Open Addressing Methods

LF	Linear Probing		Quadratic Probing		Double Hashing	
	Successful	Unsuccessful	Successful	Unsuccessful	Successful	Unsuccessful
0.05	1.0	1.1	1.0	1.1	1.0	1.1
0.10	1.1	1.1	1.1	1.1	1.1	1.1
0.15	1.1	1.2	1.1	1.2	1.1	1.2
0.20	1.1	1.3	1.1	1.3	1.1	1.2
0.25	1.2	1.4	1.2	1.4	1.2	1.3
0.30	1.2	1.5	1.2	1.5	1.2	1.4
0.35	1.3	1.7	1.3	1.6	1.2	1.5
0.40	1.3	1.9	1.3	1.8	1.3	1.7
0.45	1.4	2.2	1.4	2.0	1.3	1.8
0.50	1.5	2.5	1.4	2.2	1.4	2.0
0.55	1.6	3.0	1.5	2.5	1.5	2.2
0.60	1.8	3.6	1.6	2.8	1.5	2.5
0.65	1.9	4.6	1.7	3.3	1.6	2.9
0.70	2.2	6.1	1.9	3.8	1.7	3.3
0.75	2.5	8.5	2.0	4.6	1.8	4.0
0.80	3.0	13.0	2.2	5.8	2.0	5.0
0.85	3.8	22.7	2.5	7.7	2.2	6.7
0.90	5.5	50.5	2.9	11.4	2.6	10.0
0.95	10.5	200.5	3.5	22.0	3.2	20.0

The relative efficiency of four collision-resolution methods



Hashing Applications

- Compilers use hash tables to implement the *symbol table* (a data structure to keep track of declared variables).
- Game programs use hash tables to keep track of positions it has encountered (*transposition table*)
- Online spelling checkers.

Summary

- Hash tables can be used to implement the insert and find operations in constant average time.
 - it depends on the load factor not on the number of items in the table.
- It is important to have a prime TableSize and a correct choice of load factor and hash function.
- For separate chaining the load factor should be close to 1.
- For open addressing load factor should not exceed 0.5 unless this is completely unavoidable.
 - Rehashing can be implemented to grow (or shrink) the table.