

Unit 5 - Trie Trees & Hashing

Class notes by Vibha Masti

Feedback/corrections: vibha@pesu.pes.edu

Unit 5 - Trie Trees & Hashing

Trie Trees

Introduction

Code Implementation

Structure for Trie

Insert a pattern/key/word

Search for a pattern/key/word

Display all patterns/keys/words

Using global variables

Delete a pattern/key/word

Using global variables

Hashing & Hash Tables

Hash Function

Collision

Load factor

Collision Resolution Techniques

Separate Chaining

Code Implementation

Structure for the hash table and nodes

Initialise and destroy tables

Insert into hash table with hash function key % size

Delete from hash table

Display hash table

Open Addressing

Linear Probing

Challenges in Linear Probing

Code Implementation

Structure for hash table

Creating and Destroying the table

Insert an element

Search for an element

Delete an element

Quadratic Probing

Challenges in Linear Probing

Double Hashing

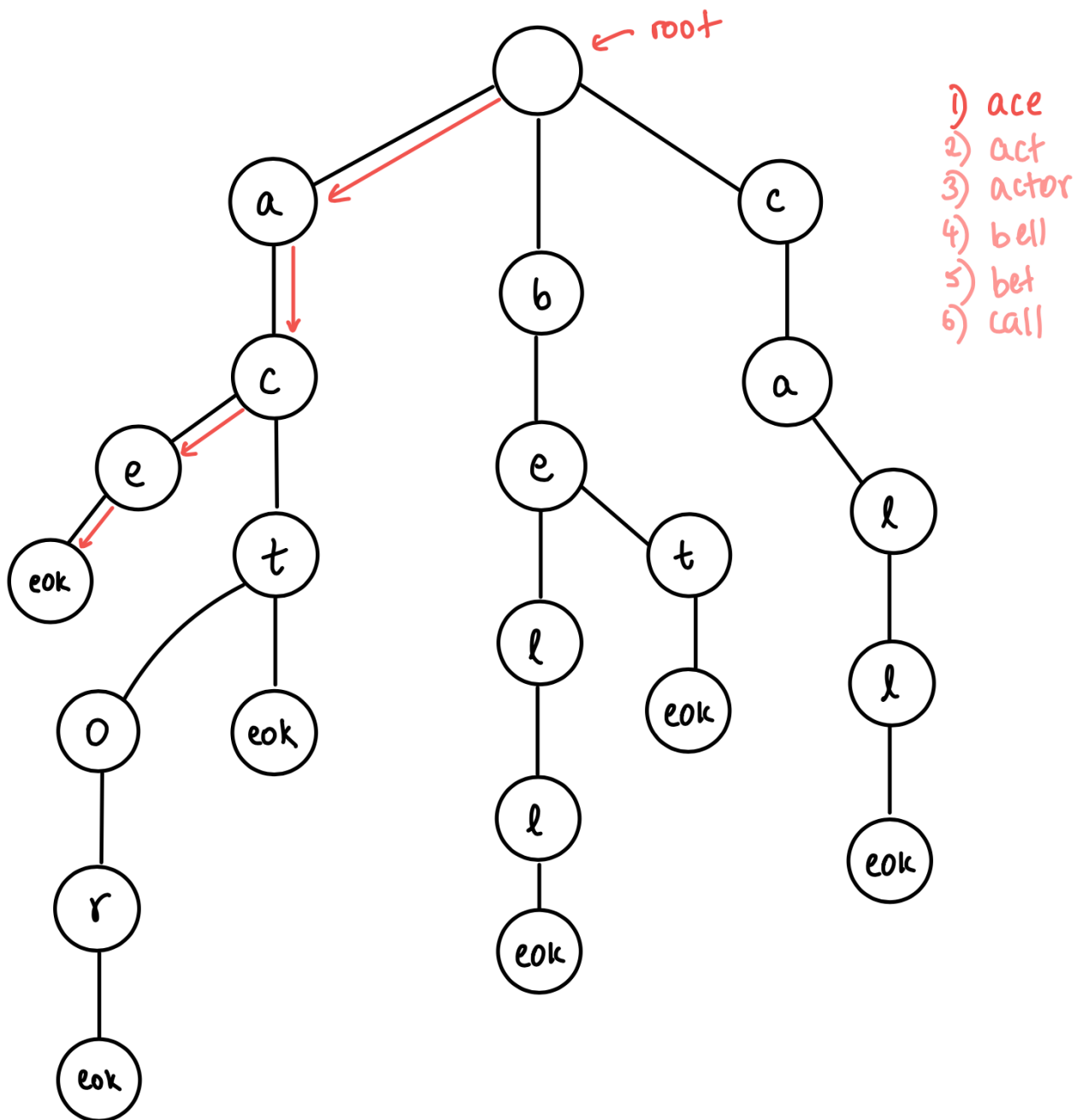
Rehashing & Load Factor

Trie Trees

- Focus on accessing data faster
- Large databases

Introduction

- Principle of BST but as many children as possible
- Comes from the word **retrieval**
- Also called prefix tree, digital tree
- Generally used to store things
- Each node can contain m pointers corresponding to m possible symbols in each position of the key
- Mobile phone databases
- If keys are strings, BSTs will compare entire strings while trie compares individual characters
- A trie is a tree where each node stores a bit indicating whether the string spelled out to this point is in the set
- Example (tracing out the word ace)



- Looking for word patterns in text

Code Implementation

- Every character in the pattern is indicated by a **link** and not a node
- If a connection exists, then the previous node exists in the pattern
- Note: I hate global variables so I have avoided using them, making the code slightly more complicated

Structure for Trie

- ASCII value of each character is the index in the `child` array of pointers
- `endofword` indicates whether that character is the end of the word in that trie pattern
- `create_node` creates a new `TrieNode` and initialises all the members of the `child` array to `NULL`

```
1  #define NO_OF_CHARACTERS 255
2  #define MAX 255
3
4  typedef struct trienode {
5      /* Child can have upto 255 children (pointers) */
6      struct trienode *child[NO_OF_CHARACTERS]; /* Array of pointers */
7      int endofword;
8  } TrieNode;
9
10
11 TrieNode *create_node() {
12     /* Create a new TrieNode */
13     TrieNode *temp;
14     temp = (TrieNode *) malloc(sizeof(TrieNode));
15     temp->endofword = 0;
16
17     /* Initialise all pointers to NULL */
18     for (int i = 0; i < NO_OF_CHARACTERS; ++i) {
19         temp->child[i] = NULL;
20     }
21     return temp;
22 }
```

Insert a pattern/key/word

```
1  void insert(TrieNode *root, char *key) {
2      /* ASCII values are indices */
3      TrieNode *cursor = root;
4      int i;
5
6      for (i = 0; key[i]; ++i) {
7          /* Find if node for that letter exists */
8          if (cursor->child[key[i]] == NULL) {
9              cursor->child[key[i]] = create_node();
10         }
11         cursor = cursor->child[key[i]];
12     }
13     cursor->endofword = 1;
14 }
```

Search for a pattern/key/word

```
1  int search(TrieNode *root, char *key) {
2      TrieNode *cursor = root;
3
4      for (int i = 0; key[i]; ++i) {
5          /* Character in that pattern absent */
6          if (cursor->child[key[i]] == NULL) {
7              /* Not found */
8              return 0;
9          }
10         cursor = cursor->child[key[i]];
11     }
12
13     /* Prefixes are not end of words */
14     return cursor->endofword;
15 }
```

Display all patterns/keys/words

```
1  void display(TrieNode *root) {
2      /* DFS traversal */
3      char word[MAX] = "";
4      int length = 0;
5      dfs_word(root, word, &length);
6  }
7
8  void dfs_word(TrieNode *root, char *word, int *plength) {
9      for (int i = 0; i < NO_OF_CHARACTERS; ++i) {
10         /* If char exists in a pattern */
11         if (root->child[i] != NULL) {
12             /* Add it to the word array */
13             word[*plength] = i;
14             ++*plength;
15             /* If end of word */
16             if (root->child[i]->endofword) {
17                 /* Print word */
18                 word[*plength] = 0;
19                 printf("%s\n", word);
20             }
21             dfs_word(root->child[i], word, plength);
22         }
23     }
24     --*plength;
```

Using global variables

```

1  /* Global variables */
2  char word[100];          // word
3  int length;              // Length of the word (search)
4  int top;                 // Top pf the stack
5  Stack s[100];           // Stack
6
7
8  void display(TrieNode *curr) {
9      for (int i = 0; i < 26; ++i) {
10         /* If char exists in a pattern */
11         if (curr->child[i] != NULL) {
12             /* Add it to the word array */
13             word[length] = ('a' + i);
14             ++length;
15             /* If end of word */
16             if (curr->child[i]->endofword) {
17                 /* Print word */
18                 for (int j = 0; j < length; ++j) {
19                     printfun(word[j]);
20                 }
21                 printfun("\n");
22             }
23             display(curr->child[i]);
24         }
25     }
26     --length;
27 }

```

Delete a pattern/key/word

- Stack needed to keep track of traversal path (all the letters in the word to be deleted)

```

1  /* struct for deletion */
2  typedef struct stack {
3      TrieNode *trie_node;    /* TrieNode in the path */
4      int index;              /* Index for removing edge */
5  } Stack;

```

- Helper functions for stack (push and pop) and check (for checking if node has connections)

```

1 void push(Stack *stack_array, TrieNode *curr, int index, int *ptop) {
2     Stack new_el;
3     new_el.index = index;
4     new_el.trie_node = curr;
5
6     stack_array[++(*ptop)] = new_el;
7 }
8
9 Stack pop(Stack *stack_array, int *ptop) {
10     return stack_array[(*ptop)--];
11 }
12
13 int check(TrieNode *curr) {
14     /* Number of connections */
15     int count = 0;
16     for (int i = 0; i < NO_OF_CHARACTERS; ++i) {
17         if (curr->child[i] != NULL) {
18             ++count;
19         }
20     }
21     return count;
22 }

```

- Function to delete

```

1 void delete_word(TrieNode *root, char *key) {
2     TrieNode *curr = root;
3     Stack stack_array[MAX], stack_el;
4     int index, top = -1, other_keys;
5
6     for (int i = 0; key[i]; ++i) {
7         /* Index of connection */
8         index = key[i];
9         if (curr->child[index] == NULL) {
10             /* No connection at that letter's index */
11             printf("word not found\n");
12             return;
13         }
14         /* Push the node and index (current letter) to stack */
15         push(stack_array, curr, index, &top);
16         curr = curr->child[index];
17     }
18     /* No longer end of word */
19     curr->endofword = 0;
20 }

```

```

21     /* Push final node to stack with letter index -1 (no letter follows)
    */
22     push(stack_array, curr, -1, &top);
23
24     /* Remove edges (connections) from nodes on stack */
25     while (1) {
26         stack_el = pop(stack_array, &top);
27
28         /* Check if not last char of word (only from second iteration) */
29         if (stack_el.index != -1) {
30             /* Make connection of node index NULL (delete last letter) */
31             stack_el.trie_node->child[stack_el.index] = NULL;
32         }
33
34         /* Root element - all chars have been deleted */
35         if (stack_el.trie_node == root) {
36             break;
37         }
38
39         /* If the trie node has other key connections */
40         other_keys = check(stack_el.trie_node);
41
42         /* If the trie node has other key
43         connections or is the end of another word */
44         if (other_keys >= 1 || stack_el.trie_node->endofword) {
45             break;
46         }
47         else {
48             /* Delete node if no connections */
49             free(stack_el.trie_node);
50         }
51     }
52 }

```

Using global variables

```

1  /* Global variables */
2  char word[100];           // word
3  int length;               // Length of the word (search)
4  int top;                  // Top pf the stack
5  Stack s[100];            // Stack
6
7
8  void delete_trie(TrieNode *root, char *key) {
9      TrieNode *curr = root;
10     int index, other_keys;
11     Stack stack_el;

```



```

12
13     top = -1;
14     for (int i = 0; key[i]; ++i) {
15         index = key[i] - 'a';
16         if (curr->child[index] == NULL) {
17             return;
18         }
19
20         push(curr, index);
21         curr = curr->child[index];
22     }
23     curr->endofword = 0;
24     push(curr, -1);
25
26     while (1) {
27         stack_el = pop();
28
29         /* Remove next connection */
30         if (stack_el.index != -1) {
31             stack_el.m->child[stack_el.index] = NULL;
32         }
33
34         /* Root element - all chars have been deleted */
35         if (stack_el.m == root) {
36             break;
37         }
38
39         other_keys = check(stack_el.m);
40
41         /* If the trie node has other key
42         connections or is the end of a word */
43         if (other_keys >= 1 || stack_el.m->endofword) {
44             break;
45         }
46         else {
47             free(stack_el.m);
48         }
49     }
50 }
51
52 int check(TrieNode *x) {
53     int count = 0;
54     for (int i = 0; i < 26; ++i) {
55         if (x->child[i] != NULL)
56             ++count;
57     }
58     return count;
59 }
60

```

```

61 void push(TrieNode *x, int k) {
62     Stack new_el;
63     new_el.index = k;
64     new_el.m = x;
65
66     s[++top] = new_el;
67 }
68
69 Stack pop() {
70     return s[top--];
71 }

```

Hashing & Hash Tables

- Retrieving data in constant time $O(1)$ instead of linear or logarithmic time
- Fast retrieval, more efficient than arrays
- Direct mapping (one-to-one) works only for numeric data of similar range, not too many gaps
- In other words, element 6 will be at 6th index of array; retrieval will be in constant time
- For more versatile data (different ranges, datatypes), direct mapping will not work as efficiently
- Too many empty locations and large ranges in such an implementation
- Need a good hash table & hash function

Hash Function

- Any function that can be used to map data of arbitrary size to fixed-size values
- Value returned by a hash function: *hash value*, *hash code*, *digest*, or simply *hash*
- For example, *mod* function, *identity* function (one-to-one), multiplicative hashing, folding and adding (adding all digits), truncation etc
- A good hash function is one that distributes keys evenly among all slots/indices (locations)

Collision

- When multiple keys generate the same hash value after being entered in to a hash function
- For example, if hash function is *mod* 10 of a number and the two numbers are 21 and 411, both result in the same hash value (1) and they collide
- Collision resolution techniques
- Hash functions that have minimum number of collisions are good
- Look up birthday paradox (24 for probability > 0.5)

Load factor

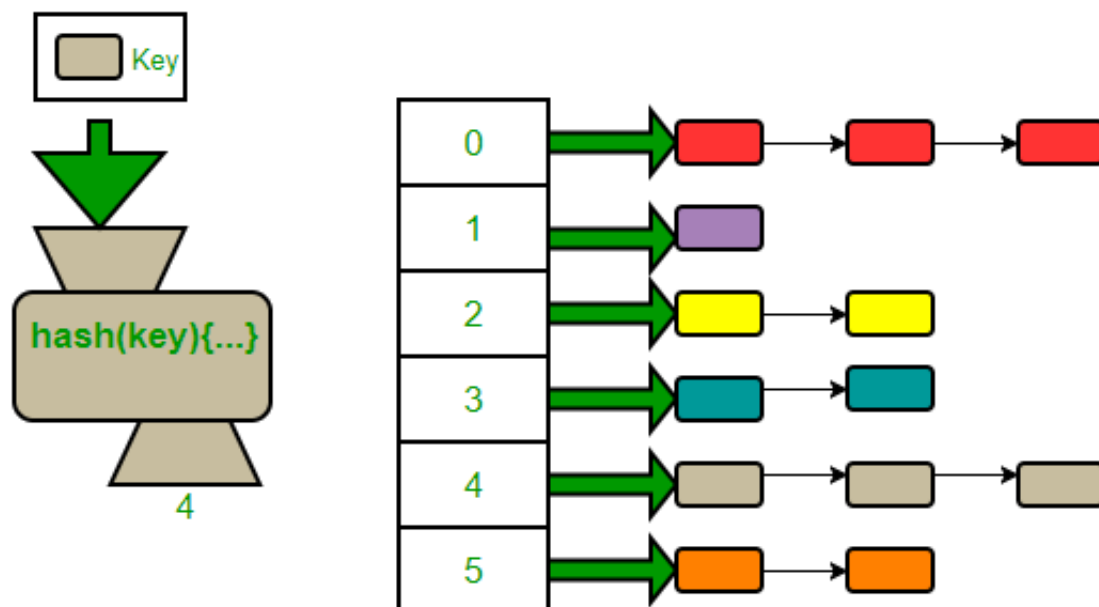
- $load\ factor = \frac{number\ of\ records\ filled}{total\ capacity}$

Collision Resolution Techniques

- Two techniques for resolving collisions
 - **Open hashing** - separate chaining (linked list at position)
 - **Closed hashing** - open addressing (all keys stored inside hash table; no new memory created)
 - Linear probing
 - Quadratic Probing
 - Double Hashing
- Still almost constant time with collisions
- Hash tables can store data in the form of key-value pairs
- For simplicity, implementation is only key
- Value can be added to the definition of the struct node

Separate Chaining

- New memory allocated during collisions



- Worst case scenario is if all the keys get mapped to the same bucket and we have a linked list of n (size of hash table) size from one single bucket, with all the other buckets empty

Code Implementation

Structure for the hash table and nodes

```
1  #define STR_LEN 20
2
3  /* Every element of the linked list */
4  typedef struct node {
5      int key;          /* int value representing the key */
6      struct node *next; /* Link */
7  } Node;
8
9  /* Every element of the hash table */
10 typedef struct hash {
11     struct node *head; /* Head of linked list */
12     int count;         /* Number of nodes in linked list */
13 } HashTable;
```

Initialise and destroy tables

```
1  HashTable *initialise_table(int size) {
2      /* Initialise array of size size */
3      HashTable *temp = calloc(size, sizeof(HashTable));
4
5      for (int i = 0; i < size; ++i) {
6          temp[i].head = NULL;
7          temp[i].count = 0;
8      }
9      return temp;
10 }
11
12 void destroy_hash(HashTable *hashtable, int size) {
13     /* Destroy hash table - delete all open chains */
14     Node *temp = NULL, *to_del = NULL;
15
16     for (int i = 0; i < size; ++i) {
17         temp = hashtable[i].head;
18         hashtable[i].head = NULL;
19
20         for (int j = 0; j < hashtable[i].count; ++j) {
21             to_del = temp;
22             temp = temp->next;
23             free(to_del);
24         }
25     }
26     free(hashtable);
```

```
27 | }
```

Insert into hash table with hash function key \% size

```
1 void insert_hash(HashTable *hashtable, int size, int key) {
2
3     /* Hash function: key % size */
4     int hash = key % size;
5
6     /* Initialise new node */
7     Node *new_node = (Node *)malloc(sizeof(Node));
8     new_node->key = key;
9
10    /* Insert to front */
11    new_node->next = hashtable[hash].head;
12
13    /* Make new node the new head */
14    hashtable[hash].head = new_node;
15
16    /* Increment count */
17    ++hashtable[hash].count;
18 }
```

Delete from hash table

```
1 void delete_hash(HashTable *hashtable, int size, int key) {
2     /* Hash function: key % size */
3     int hash = key % size;
4
5     Node *temp = hashtable[hash].head, *prev = NULL;
6
7     while (temp != NULL) {
8         if (temp->key == key) {
9
10            /* Delete node */
11            if (prev != NULL) {
12                /* Not first node */
13                prev->next = temp->next;
14                --hashtable[hash].count;
15                free(temp);
16                return;
17            }
18            else {
19                hashtable[hash].head = temp->next;
```

```

20         --hashtable[hash].count;
21         free(temp);
22         return;
23     }
24 }
25 prev = temp;
26 temp = temp->next;
27 }
28 printf("Key not found in hash table\n");
29 }

```

Display hash table

```

1 void display_hash(HashTable *hashtable, int size) {
2     /* Hash function: key % size */
3
4     Node *temp = NULL;
5     printf("\n");
6     for (int i = 0; i < size; ++i) {
7         temp = hashtable[i].head;
8         printf("%d: ", i);
9         for (int j = 0; j < hashtable[i].count; ++j, temp = temp->next) {
10             printf("%d->", temp->key);
11         }
12         printf("NULL\n");
13     }
14     printf("\n");
15 }

```

Open Addressing

- All elements are stored in the hash table itself
- In case of collision, searches for a new empty spot

Linear Probing

- Linearly probe for next empty spot
- If $hash(key) = h$ is occupied, $h + 1$ is checked and so on until $(h + i) \% n$ is found to be empty
- Example: let $hash(key) = key \% 7$ and let the elements to insert be 2, 3, 5, 7, 10, 11

$$2 \div 7 = 2 \checkmark$$

0	1	2	3	4	5	6
		2				

$$3 \div 7 = 3 \checkmark$$

0	1	2	3	4	5	6
		2	3			

$$5 \div 7 = 5 \checkmark$$

0	1	2	3	4	5	6
		2	3		5	

$$7 \div 7 = 0 \checkmark$$

0	1	2	3	4	5	6
7		2	3		5	

$$10 \div 7 = 3 \times$$

$$+1 = 4 \checkmark$$

0	1	2	3	4	5	6
7		2	3	10	5	

$$11 \div 7 = 4 \times$$

$$+1 = 5 \times$$

$$+1 = 6 \checkmark$$

0	1	2	3	4	5	6
7		2	3	10	5	11

Challenges in Linear Probing

1. **Primary Clustering:** many consecutive elements form groups and time taken to find a free slot or to search an element increases
2. **Secondary clustering:** two records have the same collision chain (Probe Sequence) if their initial position is the same

Code Implementation

Structure for hash table

```
1  typedef struct {
2      int *table;      /* int array - hash table */
3      int size;        /* size */
4  } HashTable;
```

Creating and Destroying the table

```
1  HashTable *create_table(int size) {
2      /* Initialise HashTable pointer */
3      HashTable *temp = malloc(sizeof(HashTable));
4      /* Initialise array of size size */
5      temp->table = calloc(size, sizeof(int));
6
7      /* Initialise to -1 (empty) */
8      for (int i = 0; i < size; ++i) {
9          temp->table[i] = -1;
10     }
11     temp->size = size;
12     return temp;
13 }
14
15 void destroy_table(HashTable *htable) {
16     /* Destroy hash table - free up int array */
17     htable->size = 0;
18     free(htable->table);
19 }
```

Insert an element

```
1  void insert(HashTable *htable, int element) {
2      int hash = element % htable->size;
3      int count = 0;
4
5      while (count < htable->size) {
6          if (htable->table[hash] == -1) {
7              /* Empty spot found */
8              htable->table[hash] = element;
9              break;
10         }
11
12         /* hash = hash + 1 */
13         ++hash;
14         ++count;
15     }
```



```

15
16     if (hash == htable->size) {
17         hash = 0;
18     }
19 }
20 /* Hash table is full */
21 }

```

Search for an element

```

1  int search(HashTable *htable, int element) {
2      int hash = element % htable->size;
3      int count = 0;
4
5      while (count < htable->size) {
6          /* element found */
7          if (htable->table[hash] == element) {
8              return 1;
9          }
10         ++hash;
11         ++count;
12     }
13     return 0;
14 }

```

Delete an element

```

1  void delete (HashTable *htable, int element) {
2      int hash = element % htable->size;
3      int count = 0;
4
5      while (count < htable->size) {
6          if (htable->table[hash] == element) {
7              htable->table[hash] = -1;
8              return;
9          }
10         ++hash;
11         ++count;
12     }
13 }

```

Quadratic Probing

- If $hash(key) = h$ is occupied, $h + 1^2$ is checked and so on until $(h + i^2) \% n$ is found to be empty
- The difference between consecutive squares is an odd number (the difference between consecutive numbers is 1, so every iteration we add 1 in linear probing)
- In quadratic probing, if $hash(key) = h$ is occupied, we add 1 and check, then add 3 and check, then 5, 7, 9 and so on (making sure to $\% n$). In this manner we visit all $(h + i^2) \% n$ elements
- Example: let $hash(key) = key \% 7$ and let the elements to insert be 2, 3, 5, 7, 10, 11 (same example)

$$2 \% 7 = 2 \quad \checkmark$$

0	1	2	3	4	5	6
		2				

$$3 \% 7 = 3 \quad \checkmark$$

0	1	2	3	4	5	6
		2	3			

$$5 \% 7 = 5 \quad \checkmark$$

0	1	2	3	4	5	6
		2	3		5	

$$7 \% 7 = 0 \quad \checkmark$$

0	1	2	3	4	5	6
7		2	3		5	

$$10 \% 7 = 3 \quad \times$$

$$+ 1 = 4 \quad \checkmark$$

0	1	2	3	4	5	6
7		2	3	10	5	

$$11 \% 7 = 4 \quad \times$$

$$+ 1 = 5 \quad \times$$

$$+ 3 = 8 \% 7 = 1$$

0	1	2	3	4	5	6
7	11	2	3	10	5	

Challenges in Linear Probing

1. Not all hash table slots will be visited, leaving some slots possibly empty but not checked

Double Hashing

- Applying a second hash function to key when a collision occurs
- If $hash(key) = h$ is occupied, $(h + 1 \times hash_2(key)) \% n$ is checked and so on until $(h + i \times hash_2(key)) \% n$ is found to be empty
- A popular second hash function is $hash_2(key) = PRIME - (key \% PRIME)$ where $PRIME$ is a prime number smaller than table size n

Rehashing & Load Factor

- $load\ factor = \frac{number\ of\ records\ filled}{total\ capacity}$
- When the load factor increases to more than its pre-defined value (default value of load factor is 0.75), the complexity increases
- To overcome this, the size of the array is increased (doubled) and all the values are hashed again and stored in the new double sized array to maintain a low load factor and low complexity