

16

Aim:- To find a hashed file organization strategy than can be effectively used to store and manage the movie database record.

:- To find a proper hash function

Given:- Movie database record

| IMdb ID | Title             | Year Released | gross(M) | Rating |
|---------|-------------------|---------------|----------|--------|
| 100     | Practical Magic   | 1998          | 46       | 7      |
| 101     | Dark Knight       | 2008          | 532      | 9      |
| 102     | Beetle juice      | 1998          | 73       | 7      |
| 103     | Heart breakers    | 2001          | 40       | 6      |
| 104     | Shrek             | 2001          | 267      | 8      |
| 105     | The Simpson movie | 2007          | 183      | 8      |
| 106     | The holiday       | 2006          | 63       | 7      |
| 107     | No time to die    | 2021          | 520      | 8      |
| 108     | Mr. Queen         | 2021          | 340      | 8      |
| 109     | Captain Fantastic | 2016          | 280      | 8      |
| 110     | Gifted            | 2008          | 102      | 8      |
| 111     | The family        | 2013          | 6        | 7      |
| 112     | Aladdin           | 2019          | 143      | 7      |

### Hashed file organization

Hash file organization uses the computation of hash function on some fields of the records. The hash function's output determines the location of block where the records are kept.

Extendible hashing method can be used to store and manage the database.

Static hashing is not used as database size in static hashing is fixed and this can cause certain issues. If initial no. of buckets is too small, the performance will degrade due to too much overflow. While on the other hand if very large space is allotted, there will a wastage of memory resources. Thus a dynamic hashing is preferred over it.

Linear hashing, being a dynamic hashing method could have been a probable choice. However compared with Extendible hashing, Linear hashing does not use a bucket directory and when an overflow occurs its not always the overflowed bucket that splits. In case, if no. of movies in our DB increases, there are high chances that a lot of overflow buckets get created. Although linear hashing has some advantage in certain areas like queries involving exact match but for our case extendible hashing strategy is also a good choice.

Thus, we can effectively manage and store our DB for movie records with extendible hashing strategy.

Extendible hashing uses hash function, directories and buckets to hash data and store the records in a random yet uniform way (provided a good hash function is chosen). For this case, I would be hashing movie title to allot random keys to records.

### Hash function

For this DB, polynomial rolling hash function is used. It is defined as :- For a given movie title " $s$ "  $\Rightarrow$

$$\text{hash}(s) = (s(0) + s(1) \cdot p + s(2) \cdot p^2 + \dots + s(n-1) \cdot p^{n-1}) \bmod m$$

Where  $p$  and  $m$  are some chosen positive number. It is reasonable to make  $p$  a prime number roughly equal to no. of input characters. (Usually all movie names are smaller than 31 characters (lets assume))

$\therefore \underline{p = 31}$

$m$  should be a large number since probability of two random strings colliding is  $\approx \frac{1}{m}$ . Therefore,

let  $\underline{m = 10^9 + 9}$

### Generating hash values

The hashing could be done manually but since the computation become quite comple, one could write a simple C++ code to generate the required hash values.

The code is as given below :-

```
1 // C++ Code to find hash value
2 // and binary representation of that number
3 #include <bits/stdc++.h>
4 using namespace std;
5
6 long long compute_hash(string s) {
7     const int p = 31;
8     const int m = 1e9 + 9;
9     long long hash_value = 0;
10    long long p_pow = 1;
11    for (char c : s) {
12        hash_value = (hash_value + (int)(c) * p_pow) % m;
13        p_pow = (p_pow * p) % m;
14    }
15    return hash_value;
16 }
17
18 void decToBinary(long long int n)
19 {
20     // array to store binary number
21     int binaryNum[32];
22
23     int i = 0;
24     while (n > 0) {
25
26         // storing remainder in binary array
27         binaryNum[i] = n % 2;
28         n = n / 2;
29         i++;
30     }
31
32     // printing binary array in reverse order
33     for (int j = i - 1; j >= 0; j--)
34         cout << binaryNum[j];
35 }
36
37 int main()
38 {
39     int n;
40     cin >> n;
41     cin >> ws;
42     for (int i = 0; i < n; i++) {
43         string s;
44         getline(std::cin, s);
45
46         long long int x = compute_hash(s);
47         cout << s << " " << x << endl;
48         decToBinary(x);
49         cout << endl << endl;
50 }
```

Sat Nov 27 05:02

/home/tanmay/Desktop/abc.cpp - CP Editor

File Edit Actions View Options Help

abc.cpp x Untitled-1 x

```
1 // C++ Code to find hash value
2 // and binary representation of that number
3 #include <bits/stdc++.h>
4 using namespace std;
5
6 long long compute_hash(string s) {
7     const int p = 31;
8     const int m = 1e9 + 9;
9     long long hash_value = 0;
10    long long p_pow = 1;
11    for (char c : s) {
12        hash_value = (hash_value + (int)(c) * p_pow) % m;
13        p_pow = (p_pow * p) % m;
14    }
15    return hash_value;
16 }
17
18 void decToBinary(long long int n)
19 {
20     // array to store binary number
21     int binaryNum[32];
22
23     int i = 0;
24     while (n > 0) {
25
26         // storing remainder in binary array
27         binaryNum[i] = n % 2;
28         n = n / 2;
29         i++;
30     }
31
32     // printing binary array in reverse order
33     for (int j = i - 1; j >= 0; j--)
34         cout << binaryNum[j];
35 }
36
37 int main()
38 {
39     int n;
40     cin >> n;
41     cin >> s;
42     for (int i = 0; i < n; i++) {
43         string s;
44         getline(std::cin, s);
45
46         long long int x = compute_hash(s);
47         cout << s << " " << x << endl;
48         decToBinary(x);
49         cout << endl << endl;
50 }
```

Messages

Clear C++

[05:01:46] [Compiler] [Compilation has started]  
[05:01:47] [Compiler] [Compilation has finished]  
[05:01:47] [Runner[1]] [Execution has started]  
[05:01:47] [Runner[1]] [Execution for test case #1 has finished in 12ms]

Test Cases 0 / 0 / 1 Add Test More

Checker: Ignore trailing spaces Add Checker

| ✓ Input #1  | Run | Output #1  | ** | Expected #1 | Del |
|---|-----|--|----|-------------|-----|
| 13<br>Practical Magic<br>Dark Knight<br>Beetle Juice<br>Heart Breakers<br>Shrek<br>The Simpsons Movie<br>The Holiday<br>No Time to Die<br>Mr. Queen<br>Captain Fantastic<br>Gifted<br>The Family<br>Aladdin |     | Practical Magic 135733021<br>1000000101110001111100011101<br><br>Dark Knight 520000018<br>111101111110101100010101010<br><br>Beetle Juice 2144590<br>1000001011100101001110<br><br>Heart Breakers 491083995<br>1110101000101010100011011011<br><br>Shrek 101938499<br>110000100110111010101000011<br><br>The Simpsons Movie 565992773<br>100001101111000101110101000101<br><br>The Holiday 181816787<br>1010110101100100110111010011<br><br>No Time to Die 494532647<br>11101011110011111100000100111<br><br>Mr. Queen 850562298<br>110010101100101000110011111010<br><br>Captain Fantastic 155431945<br>1001010000111011010000001001<br><br>Gifted 359865338<br>10101011100110001101111111010 |    |             |     |

Line 50, Column 6

Compile Run Compile and Run

Notes:-  $\text{hash}(s) = (s(0) + s(1) \cdot p + \dots + s(n-1) \cdot p^{n-1}) \bmod m$ ,  
 here  $s(i)$  denotes ascii value of  $i^{\text{th}}$  character of string  
 For ex  $\Rightarrow$  ascii value of "1" is 49 and of "a" is 97 etc.

### Hash values

| <u>Title</u>       | <u>hashed no.</u> | * (only last 6 digits) (lsb)<br><u>Binary</u> |
|--------------------|-------------------|---|
| Practical magic    | 135733021         | 011101  |
| Dark Knight        | 520008018         | 010010  |
| Bee the juice      | 21445901          | 001110  |
| Heart breakers     | 491083995         | 011011  |
| Shrek              | 101938499         | 000011  |
| The Simpson movies | 565992773         | 000101  |
| The holiday        | 181816787         | 010011  |
| No time to die     | 494532647         | 100111  |
| Mr. Queen          | 850562298         | 111010  |
| Captain fantastic  | 155431945         | 001001  |
| Gifted             | 359865338         | 111010  |
| The family         | 836279542         | 110110  |
| Aladdin            | 726892695         | 010111  |

## Schematic of file organization

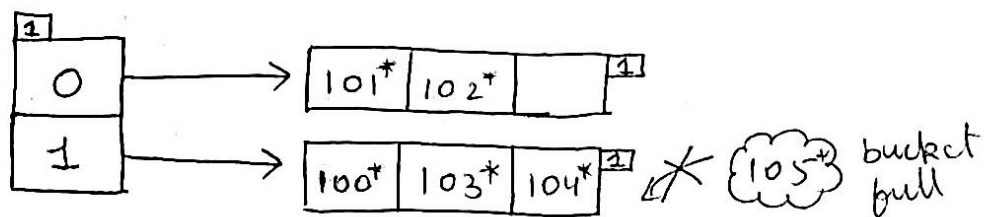
Note:- using  $(ID_i)^*$  to represent a full record i.e.  $\Rightarrow$   
 $(ID_i)^* \Rightarrow (ImdbID)_i, (Title)_i, (year\ released)_i, (gross(M))_i, (rating)_i$

### Terms used:-

1. Directories  $\Rightarrow$  These containers store pointers to bucket. No. of directories =  $2^{\text{global depth}}$ .
2. Buckets  $\Rightarrow$  They store the hashed keys
3. global depth  $\Rightarrow$  Number of bits used to categorize the key.  
global depth = no. of bits in directory.
4. local depth  $\Rightarrow$  Same as global depth but it is used for buckets.
5. overflow  $\Rightarrow$  It is the situation when the bucket can no longer hold extra records and a new record is pushed in that bucket.

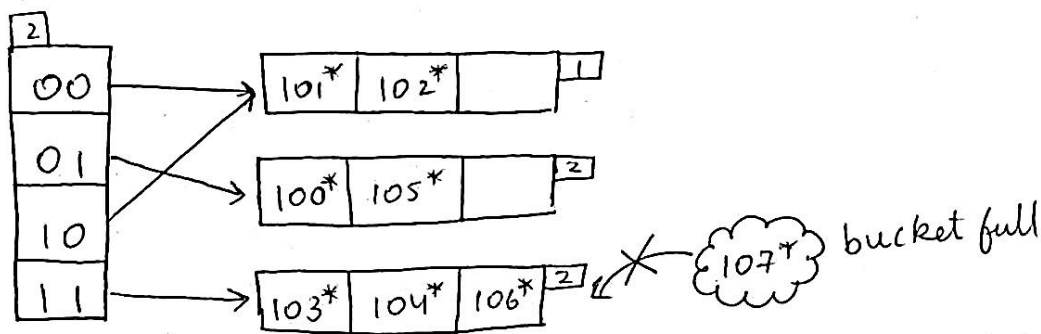
\*(Inserting data in the order of Imdb ID)

1.) Global depth = 1 \*(Note: A bucket can hold upto 3 values)



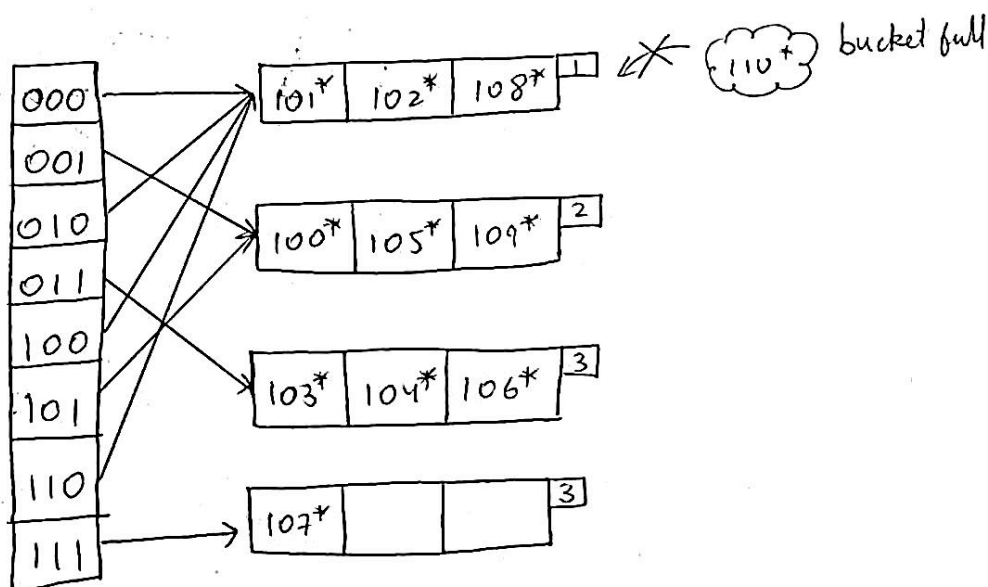
overflow on inserting  $105^*$ , since bucket is already full.  
local depth = global depth splitting occurs + Directory expansion.

global depth = 2

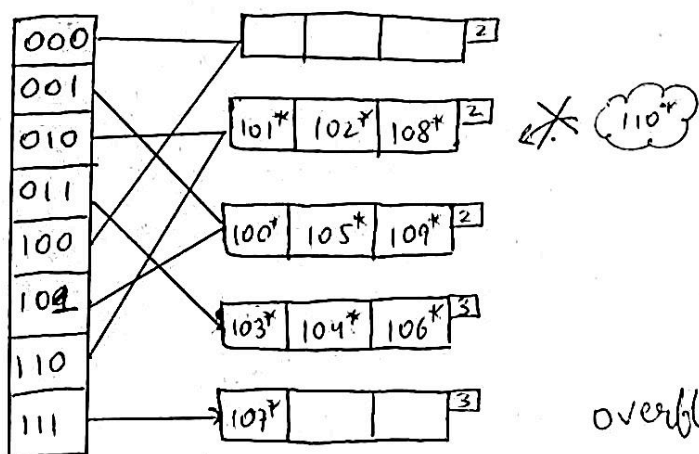


overflow on inserting  $107^*$   $\Rightarrow$  local depth = global depth  
directory expansion plus bucket splitting

global depth = 3

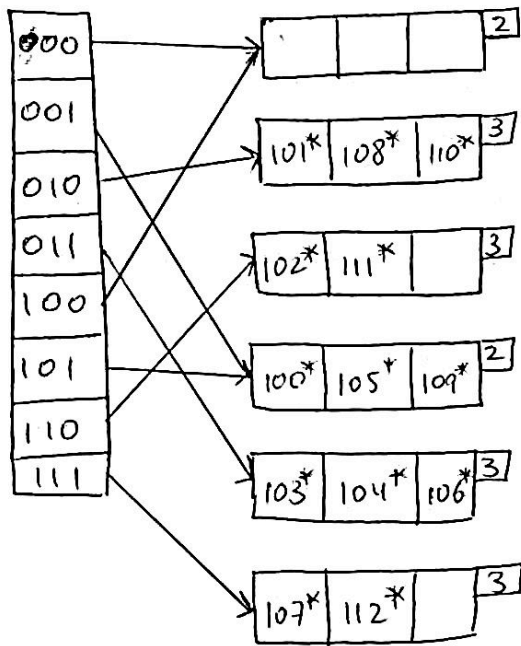


overflow on inserting  $110^*$   $\Rightarrow$  local depth < global depth  
only bucket splitting



overflow again  $\rightarrow$  local depth < global depth  
bucket split





This is the final schematic diagram of file organization for given movie database record.

Adding 5 new records

| Jmdb ID | Title           | Year Released | gross(M) | rating | hash value | Binary (61sb) |
|---------|-----------------|---------------|----------|--------|------------|---------------|
| 113     | Smurfs          | 2005          | 180      | 7      | 390063585  | 100001        |
| 114     | captain america | 2004          | 450      | 8      | 384424966  | 000110        |
| 115     | laigan          | 2001          | 300      | 9      | 104579135  | 111111        |
| 116     | border          | 2002          | 150      | 9      | 360091001  | 111 001       |
| 117     | shang-chi       | 2021          | 500      | 9      | 119638756  | 100 100       |

(P.T.O)

Activities CP Editor Sat Nov 27 05:01 /home/tanmay/Desktop/abc.cpp - CP Editor

File Edit Actions View Options Help

abc.cpp \* X Untitled-1 \* X

```
1 // C++ Code to find hash value
2 // and binary representation f that number
3 #include <bits/stdc++.h>
4 using namespace std;
5
6 long long compute_hash(string s) {
7     const int p = 31;
8     const int m = 1e9 + 9;
9     long long hash_value = 0;
10    long long p_pow = 1;
11    for (char c : s) {
12        hash_value = (hash_value + (int)(c) * p_pow) % m;
13        p_pow = (p_pow * p) % m;
14    }
15    return hash_value;
16 }
17
18 void decToBinary(long long int n)
19 {
20     // array to store binary number
21     int binaryNum[32];
22
23     int i = 0;
24     while (n > 0) {
25
26         // storing remainder in binary array
27         binaryNum[i] = n % 2;
28         n = n / 2;
29         i++;
30     }
31
32     // printing binary array in reverse order
33     for (int j = i - 1; j >= 0; j--)
34         cout << binaryNum[j];
35 }
36
37 int main()
38 {
39     int n;
40     cin >> n;
41     cin >> ws;
42     for (int i = 0; i < n; i++) {
43         string s;
44         getline(std::cin, s);
45
46         long long int x = compute_hash(s);
47         cout << x << endl;
48         decToBinary(x);
49         cout << endl << endl;
50 }
```

Messages

Clear C++

[04:40:40] [Compiler] [Compilation has started]  
[04:40:41] [Compiler] [Compilation has finished]  
[04:40:41] [Runner[1]] [Execution has started]  
[04:40:41] [Runner[1]] [Execution for test case #1 has finished in 11ms]

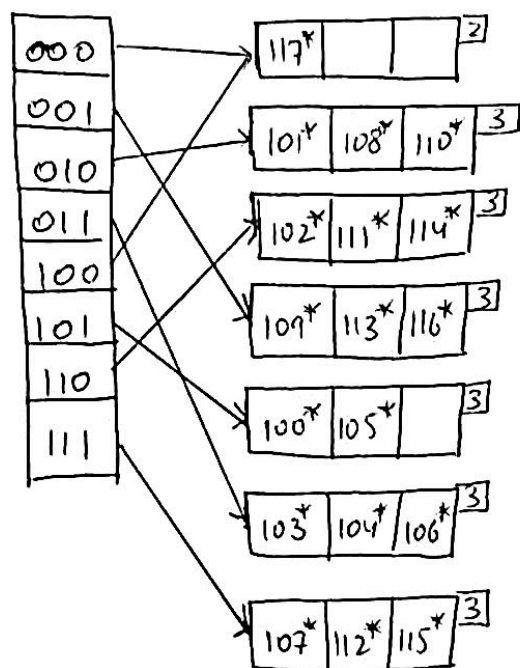
Test Cases 0 / 0 / 1 Add Test More

Checker: Ignore trailing spaces Add Checker

| ✓ Input #1   | Run | Output #1   | ** | Expected #1 | Del |
|--|-----|---|----|-------------|-----|
| 5<br>smurfs<br>captain america<br>lagan<br>border<br>shang-chi |     | smurfs 390063585<br>10111001111111110010111100001<br><br>captain america 384424966<br>10110111010011101110000000110<br><br>lagan 104579135<br>11000111011110000000011111<br><br>border 360091001<br>10101011101101000110101111001<br><br>shang-chi 119638756<br>111001000011000101011100100 |    |             |     |

Line 50, Column 6 Compile Run Compile and Run

global depth=3



This is the updated schematic diagram after 5 new records are inserted

(Note  $\Rightarrow (ImdbID)_i^*$  denotes full record containing  
 $(ImdbID)_i, (Title)_i, (year\ released)_i, (gross(M))_i, (rating)_i$ )

After adding 5 new records the hash function has generated quite random and uniform hash-values

Thus extendible hashing along with our polynomial rolling hash function can be efficiently use to store and manage this movie database.

— The End —