

## END-SEM PART-B

Que 16:-

IMDB ID	Title	Year released	gross (M)	Rating
100	Practical Magic	1998	46	7
101	Dark knight	2008	532	9
102	Beetle Juice	1998	13	7
103	Heart Breakers	2001	40	6
104	Shrek	2001	267	8
105	The Simpsons Movie	2007	183	8
106	The Holiday	2006	63	7
107	No Time To Die	2021	520	8
108	My 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

We have to find a suitable hashing strategy for the movie records:

→ Hashed File Organisation: Hash file Organisation 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.

We have three options to store this data i.e. ~~static~~ hashing, linear hashing or extendible hashing.

For our purpose we are going to use the extendible hashing and so we will see the reasons why we discarded the first two methods.

At first we cannot use static hashing because in it the database size has to be fixed which is not the case for us. In it if the initial no. of buckets is very small, then the performance will degrade due to overflow and if the number becomes too large memory allocated will be very large. Therefore we cannot use static hashing and we have to choose a method from dynamic hashing.

In Dynamic Hashing we have two options either to use linear hashing or go with extendible hashing. We will be using extendible hashing because in linear hashing -

Firstly it does not use a bucket directory.

When an overflow occurs, the pointer bucket splits which means that the bucket which has overflowed may or may not be the one to split. which means that a lot of overflow buckets are created.

So, if we see in the actual world, the no. of movies are very high, which means that if we would go with linear hashing then ~~over~~ there are a lot of chances that our buckets no. is way too high.

So, for our case we are using Extendible Hashing

Extendible Hashing

Extendible Hashing uses hash function, directories and buckets to hash data and store the records in a random yet uniform way given that we use a good hash function.

For our case we have taken the hashing of movie title to allot random keys to records.

## HASH FUNCTION

For our usecase, we will use Polynomial Rolling Hash Function

The polynomial hash function is defined as :-

lets say that the movie title is 's'

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

where,  $m, p$  are both positive numbers.

~~pr~~  $p = 31$  is generally a prime no and we choose it approximately equal to the number of input characters.

So, now we are assuming that  $p = 29$ , as mostly movie names have less than 31 characters.

$m = m$  ~~is~~ over here is a very large number since

the probability of two random strings colliding is  $\approx \frac{1}{m}$

So we are assuming  $m$  to be  $10^9 + 9$

## Generating Hash Values

Generally we do the hashing by hand, but in our case the records are quite complex and they require a lot of calculation therefore we have used a C++ code to ~~use~~ generate our hash values for the given data.

In our function

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

$s(i)$  denotes the ASCII value of the  $i^{\text{th}}$  character of string.

eg

ASCII value of	0 $\rightarrow$ 48	a $\rightarrow$ 97
	1 $\rightarrow$ 49	b $\rightarrow$ 98
	2 $\rightarrow$ 50	c $\rightarrow$ 99

The C++ code for our usecase is given below:

```
1  #include <bits/stdc++.h>
2  using namespace std;
3
4  long long findHashValue(string movie) {
5      const int m = 1e9 + 9;
6      const int p = 29;
7      long long ans = 0;
8      long long powerUp = 1;
9      for (char c : movie) {
10         ans = (ans + (int)(c) * powerUp);
11         ans = ans % m;
12         powerUp = powerUp * p;
13         powerUp = powerUp % m;
14     }
15     return ans;
16 }
17
18 void convertBinary(long long int number){
19     int representation[32];
20
21     int count = 0;
22     while (number > 0){
23         representation[count] = number % 2;
24         number = number / 2;
25         count++;
26     }
27     for (int i = count - 1; i >= 0; i--)
28         cout << representation[i];
29 }
30
31 int main(){
32     int testCase;
33     cin >> testCase;
34     cin >> ws;
35     for(int i=0; i<testCase; i++){
36         string s;
37         getline(std::cin, s);
38
39         long long int x = findHashValue(s);
40         cout << s << endl;
41         cout << x << endl;
42         convertBinary(x);
43         cout << endl << endl;
44     }
45
46     return 0;
47 }
```



Practical Magic

238009119

1110001011111011101100011111

Dark Knight

402826310

1100000000101010010001000110

Beetle Juice

787552802

101110111100010001101000100010

Heart Breakers

221877205

1101001110011001001111010101

Shrek

78241329

100101010011101111000110001

The Simpsons Movie

979252036

111010010111100011001101000100

The Holiday

365459057

10101110010000111011001110001

No Time to Die

42466832

10100001111111111000010000

Mr. Queen

877493572

110100010011010111110101000100

Captain Fantastic

730724320

101011100011011111011111100000

Gifted

159814386

1001100001101001001011110010

The Family

716535315

101010101101010111011000010011

## Computed Hash Values

Movie Title is hashed for our case and we have find the value corresponding to that. Now, as the binary representation is quite large and we here have to work for small no. of records. therefore we are using the last 5 binary digits from the representation.

ID	Title	Hashed No	Binary Representation
100 <sup>a</sup>	Practical Magic	238009119	11111
101 <sup>a</sup>	Dark Knight	402826310	00110
102 <sup>a</sup>	Beetle Juice	787552802	00010
103 <sup>a</sup>	Heart Breakers	221877205	10101
104 <sup>a</sup>	Shrek	78241329	<del>10001</del> 10001
105 <sup>a</sup>	The Simpsons Movie	979252036	00100
106 <sup>a</sup>	The Holiday	365459057	10001
107 <sup>a</sup>	No Time To Die	42466832	10000
108 <sup>a</sup>	Mr Queen	877493572	00100
109 <sup>a</sup>	Captain Fantastic	730724320	00000
110 <sup>a</sup>	Cyfrid	159814386	10010
111 <sup>a</sup>	The Family	716535315	10011
112 <sup>a</sup>	Aladdin	657487126	10110

## SCHEMATIC OF THE FILE ORGANIZATION:

Now as the entries are too big, I will be using the IMDB ID to represent one full entry i.e.

$$(ID)_i = (IMDB\ ID)_i, (Title)_i, (Year\ Released)_i, (Gross(M))_i, (Rating)_i$$

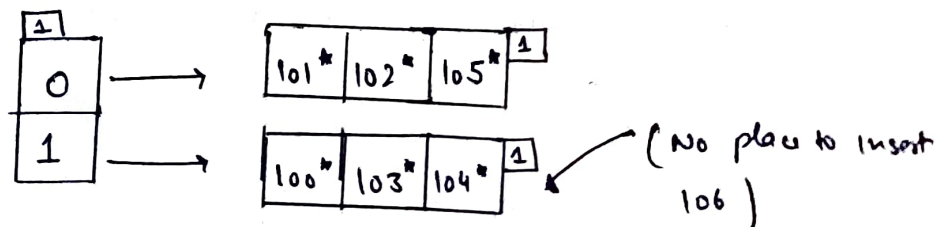
## Terms Used

1. Buckets → They are used to store the hashed keys
2. Directories → These are the containers that store pointers to bucket.  
 $\text{No. of directories} = 2^{\text{global depth}}$
3. Global Depth → It is associated with the directories. They denote the number of bits which are used by the hash function to categorize the keys.  
 $\text{Global Depth} = \text{No. of bits in directory}$
4. Local Depth → It is same as that of Global Depth except for the fact that local depth is associated with the buckets and not the directories.
5. Bucket splitting → when the number of elements in a bucket exceeds a particular size, then the bucket is split into two parts.
6. Directory Expansion → Directory Expansion takes place when a bucket overflows. It is performed when the local depth of the overflowing bucket is equal to the global depth.

## INSERTING DATA

We will be inserting data in the order of IMDB Rating.  
And we are taking the bucket capacity to be 3

1. Global Depth = 1

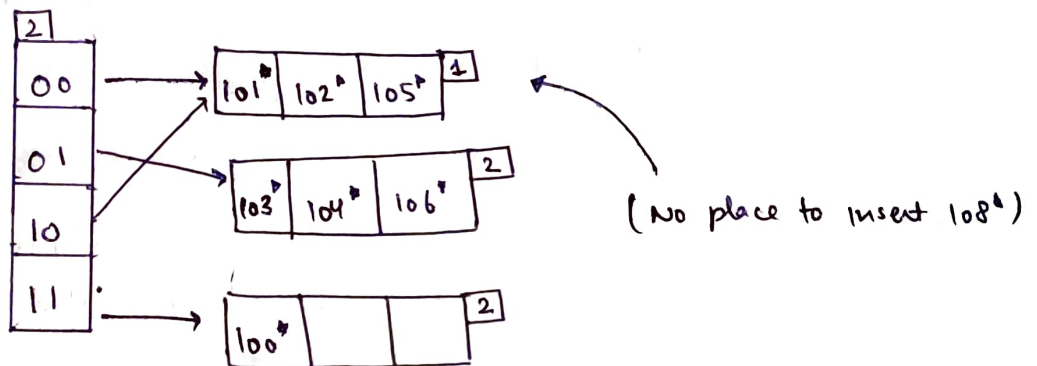


Here when we are trying to insert  $106^*$  the 0 bucket is already full so we reached an overflow.

As  
Local depth = Global Depth.

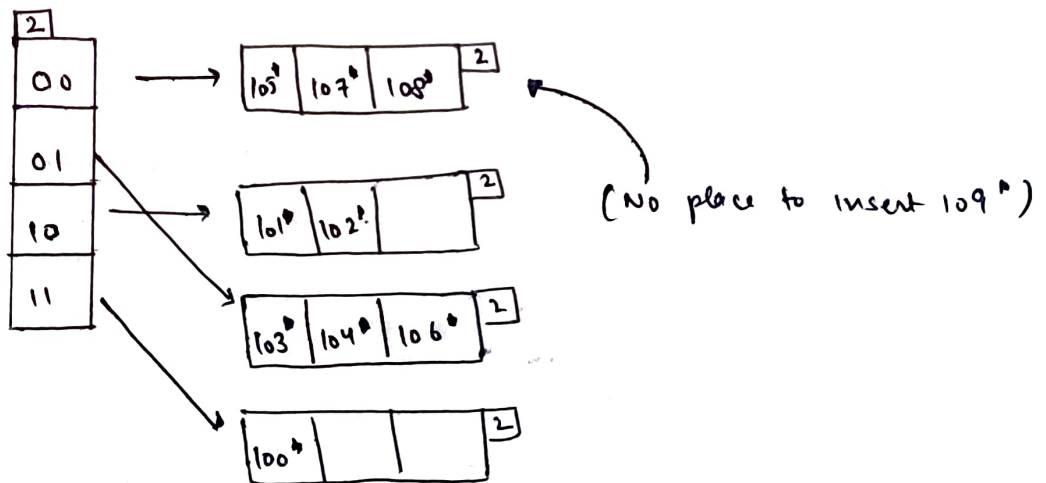
Directory expansion + splitting will happen.

Global Depth = 2



As local depth < global depth; only bucket splitting takes place.

Global Depth = 2

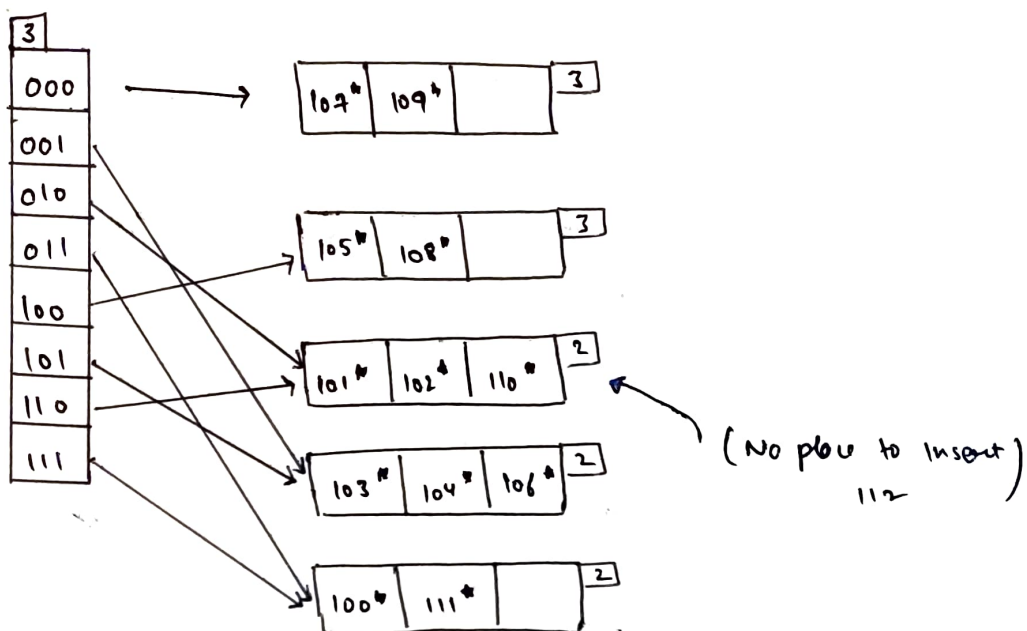


Local depth = global depth.

Directory expansion + splitting takes place.

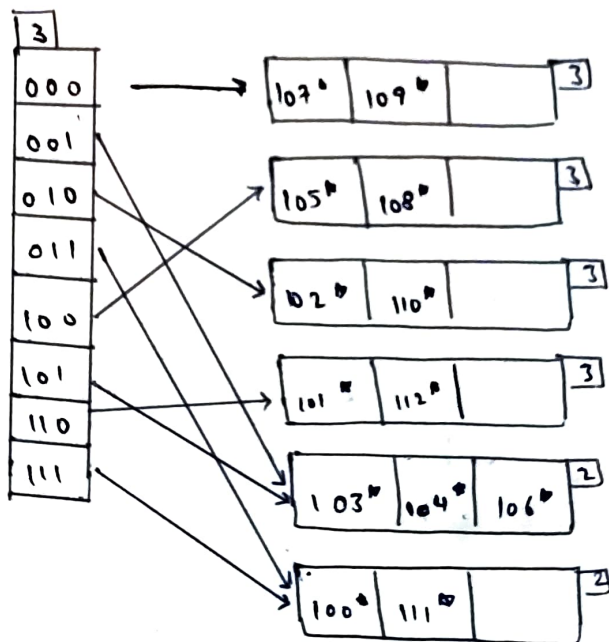


Global depth = 3



Local depth < global depth; only bucket splitting.

Global Depth = 3



This represent our final schematic diagram of file organisation for the given data.

Inserting 5 new Records

Names of Movies Chosen → Here Pheri, Golmaal, Ek Tha Tiger, Tannu Weds Manu, Kabir Singh.

Hera Pheri

169427156

1010000110010100000011010100

Golmaal

301858303

1000111111011111110111111111

Ek tha Tiger

938534474

110111111100001110011001001010

Tanu Weds Manu

213311373

1100101101101101111110001101

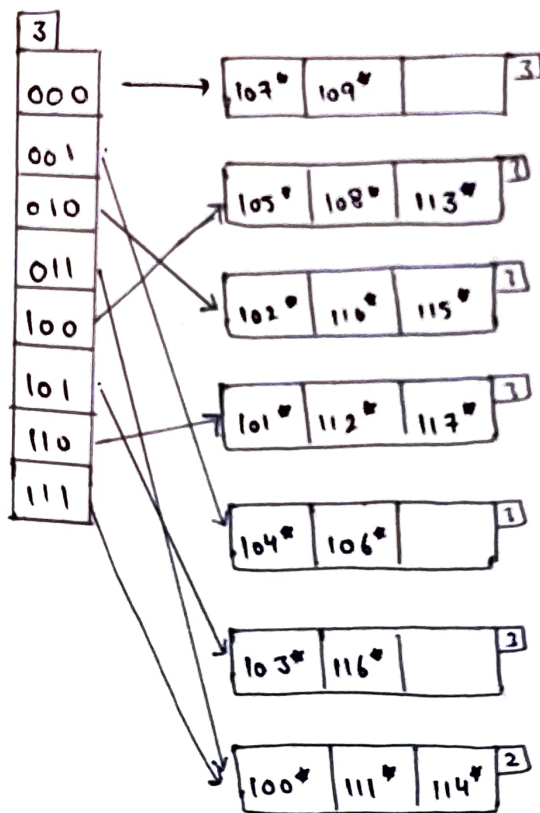
Kabir Singh

343400078

101000111011111101111010001110

IMDB ID	Title	Year Released	Questions	Rating	Hashed No.	Binary Representation
113	Herz Phoenix	2004	500	9	169427156	10100
114	Goldmal	2001	415	8	301858303	11111
115	Ek the Tiger	2012	312	9	938534474	01010
116	Tanu Weds Manu	2011	240	7	213311373	01101
117	Kabir Singh	2019	300	8	343400078	01110

Now when we will make the updated Schems. at entry 116 the bucket will split and as local depth < global depth only bucket splitting takes place.



So after 5 data entries our hash function generated some new values and the updated schematic is shown above.