

Algorithms for Big Data
Prof. John Ebenezer Augustine
Department of Computer Science and Engineering
Indian Institute of Technology, Madras

Lecture – 30
Perfect Hashing

In this segment, we are going to talk about Perfect Hashing, Using 2-universal families of hash functions.

(Refer Slide Time: 00:22)

Motivation:

S
 m elements, n -bin hash table, chain hashing
 h chosen UAR from 2-universal family

What is the expected search time?

Suppose: search item $x \notin S$ ↳ Denoted by X

$$E[X] = \sum_{i=1}^m E[X_i] \leq \sum_{i=1}^m \frac{1}{n} = \frac{m}{n} = 1 \text{ if } m=n$$

↳ indicates collision of x with item i

Let us begin some motivation. Let us say we have m elements and n bin hash table. Let us say we are performing chain hashing and let us assume the hash function h is chosen uniformly at random from at two universal family fashions. The question is what is the expected search time? Let us assume that. Let us start with the case that X is not in the dictionary. So, S is the set of items in the dictionary and let assume what else even happens when X is not in the dictionary.

The search item, let us denote the search time by the random variable X . E of X is going to be the summation over 1 to n , E_i is an indicator random variable that denotes the collision of X , where the search time X with the item i in the dictionary and just to be clear, the summation must go from i equal to 1 to M and this summation is at most summation i equal to 1 to m 1 over n because we have assumed that we have drawn from the random from a 2-universal family of hash functions. That is at most m over n and

when we limit ourselves to the case where m equals n , then the expected search time is only 1.

(Refer Slide Time: 02:07)

Motivation:

Conclusion: Expected search $O(1)$ if $m = n$

S
 m elements, n -bin hash table, chain hashing,
 h chosen UAR from 2-universal family

What is the expected search time?

Suppose: search item $x \in S$ \hookrightarrow Denoted by X

$E[X] = 1 + \sum_i E[X_i] = 1 + \frac{m-1}{n} \approx 2$ if $m=n$

\hookrightarrow Colliding with self.
 \hookrightarrow Indicates collision of x with item i
 $2 \leq i \leq m$

NPTEL

Now, what about the case where the search item X is not an s , where E can similarly bound the expectation on X , but now we need to be very careful because X is in the dictionary. We cannot sum, I mean we cannot use, indicate a random variable for the item X in the dictionary because X will collide with itself with probability 1.

So, we are going to isolate that by a separate 1 and for the remaining cases, we are going to use X_i , but this really does not affect us much because when we apply the same way of evaluating E of X_i , we get 1 plus $m-1$ over m which is approximately 2 when m equal to n and important conclusion we can make out of this exercise is that the expected search time is only O of 1 constant when m equal to n and the idea of perfect hashing is motivated from this context, but one of be a little bit more at our expectation instead of requiring an expected search time of O of 1.

(Refer Slide Time: 03:47)

Can $O(1)$ Search time be guaranteed?

Focus: Static dictionaries

Fixed set of items
No updates

Deterministically!

NPTEL 37

We are interested in a deterministic search time of $O(1)$ can achieve this requirement and we are going to live it our focus to static dictionaries in which we are going to store a fixed set of items. This is going to be no updates to these items and for this fixed set, we will not be able to query whether an item x in the universe is present in the dictionary or not.

(Refer Slide Time: 04:19)

Can $O(1)$ Search time be guaranteed?

Focus: Static dictionaries

Deterministically!

Recall: $\Pr(\text{most loaded bin} \geq m\sqrt{\frac{2}{n}}) \leq \frac{1}{2}$

(m balls into n bins using $h \in 2$ -universal family of hash functions)

We want most loaded bin + have at most 1 ball

$\Rightarrow m\sqrt{\frac{2}{n}} = 1 \Rightarrow n = 2m^2$

Therefore, if $n = 2m^2$, $\exists h$ s.t. search time is $O(1)$; Can be found in exp 2 trials.

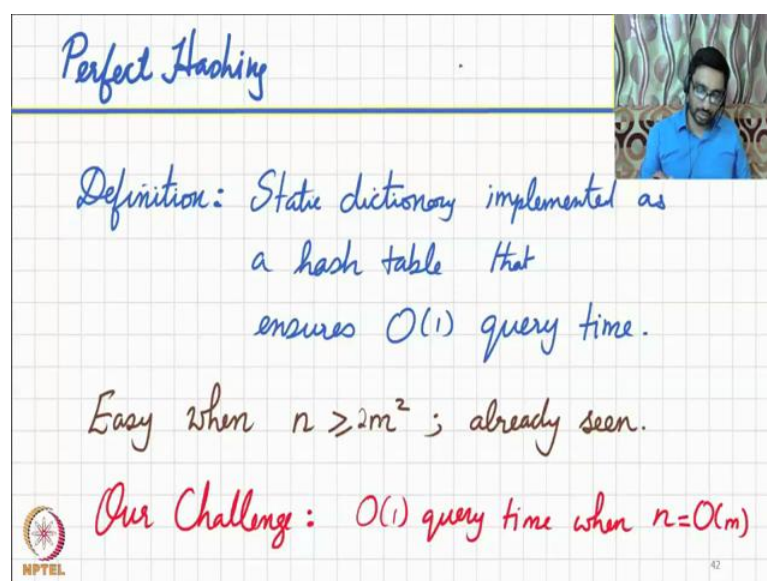
The important thing is that we are not even going to require complete randomness. We are going to assume that the hash function we have available shows uniformly at random

from a two universal family of hash functions. Before we see let us recall the balls and bins result that we have. We have m balls and n bins and these balls are hash into bins using the hash function that shows uniformly random from a family of hash functions. Then, we know that the probability that the most loaded bin is more than m times square root of 2 over n is at most a half.

So, if we want constant query time, we have to ensure that this m times 2 square root of 2 over n is 1 or lesser and this will mean that n has to be at least two times m square and when we set n equal to $2 m$ square, the probability with which we get a collision with which we get a search time of 1 is at most a half.

Suppose, we pick a hash function h from uniformly at random from this family of hash function and we hash each item in the dictionary to location using this hash function with probability at least half where going to get no collisions at all which means that if we tried to construct such static dictionary with probability half, we will succeed in getting a perfect hashing where in the search time is at most O of n and this is simple. But the trial if we want to insist on success, we simply repeat and we get success. So, with expected two times we will be able to get a static dictionary in which n items are mapped into $2 m$ square root locations and the query item for any item is at most 1.

(Refer Slide Time: 07:03)



Perfect Hashing

Definition: Static dictionary implemented as a hash table that ensures $O(1)$ query time.

Easy when $n \geq 2m^2$; already seen.

Our Challenge: $O(1)$ query time when $n = O(m)$

NPTEL

42

So, all of this discussion motivates what we call perfect hashing. So, we want to store a static dictionary and we want to have that implemented as a hash table and the key

requirement is that our query time in the words case must be $O(1)$ and as we have just seen if we are allowed to use a lot of ha, i mean much larger tables. So, if we are allowed to use a table of length at least $2m^2$, then this problem is already solved and that is easy. You just repeatedly attempt constructing the static dictionary using the straight forward hashing technique until we get this requirement.

Our question now is a little bit more challenging. How do we construct such a hash table using only little space that is proportional to the number of items and still want query time to be constant not with high probability, but deterministically.

(Refer Slide Time: 08:15)

Perfect Hashing

Remember, we only need a table of size $n=m$ to ensure that number of collisions is $\leq m$ with probability at least $1/2$

Recall: $X = \#$ of collisions ; $E[X] = \frac{m^2}{2n}$

Therefore $\Pr[X \geq \frac{m^2}{n}] \leq \Pr[X \geq 2E[X]] \leq \frac{1}{2}$
(Markov's inequality)

\Rightarrow When $n=m$, $\Pr[X \geq m] \leq \frac{1}{2}$

Algorithm: Try hash functions until # of collisions $\leq m$
Hash each chain of length c_i into a hash table of size $2c_i^2$

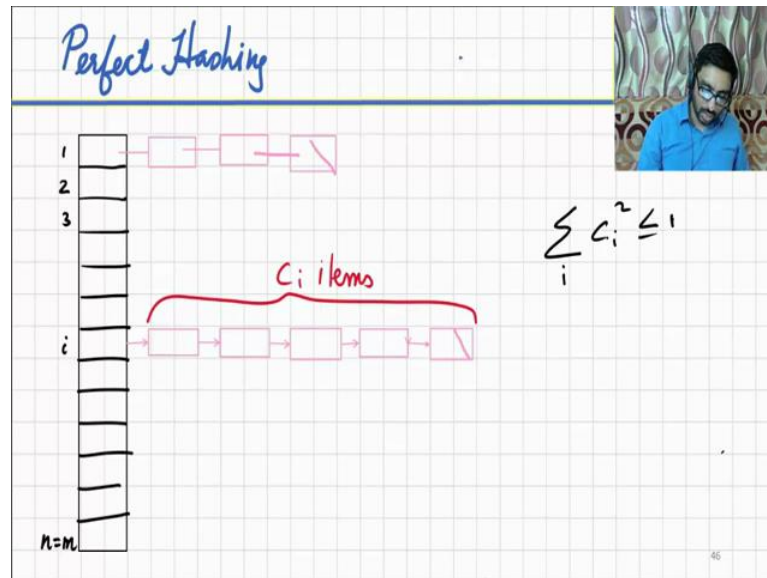
No collisions

Let us remind ourselves in the previous segment that the number of collision exceeds m with probability at most a half and when n equal to m , the number of collisions will exceed m with probability at most a half. So, here is an idea for how to construct this perfect hashing. We try hash functions until the number of collisions is at most m and of course, we can by the probabilistic method, we can keep repeating our choice of h until we find the hash function h that satisfies the requirement, but this is not sufficient for us because now there are collisions. So, some of the collisions can happen on the same bin which means that the query time need not be $O(1)$ anymore, but here is what we do.

Whenever an item collide with another item, that means that particular hash location will have multiple items in it, so let say that we create a chain of length c_i at that location i .

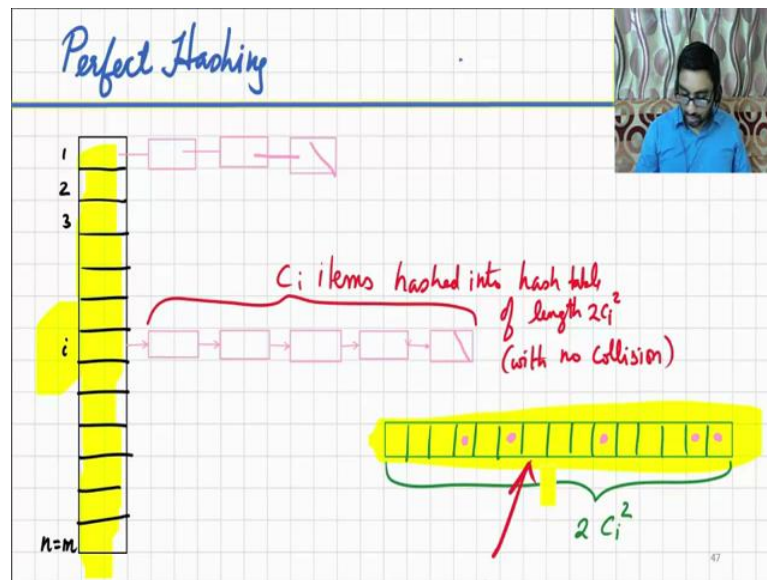
What do we do is, we take those c_i items and hash them using a different hash function into a hash table of size two c_i times c_i square and we choose the second hash function for the c_i items in such a way that there are no collisions in the second hash table, but keep in mind the second hash table has to be constructed for each location where there is multiple items that are hash into the location.

(Refer Slide Time: 10:25)



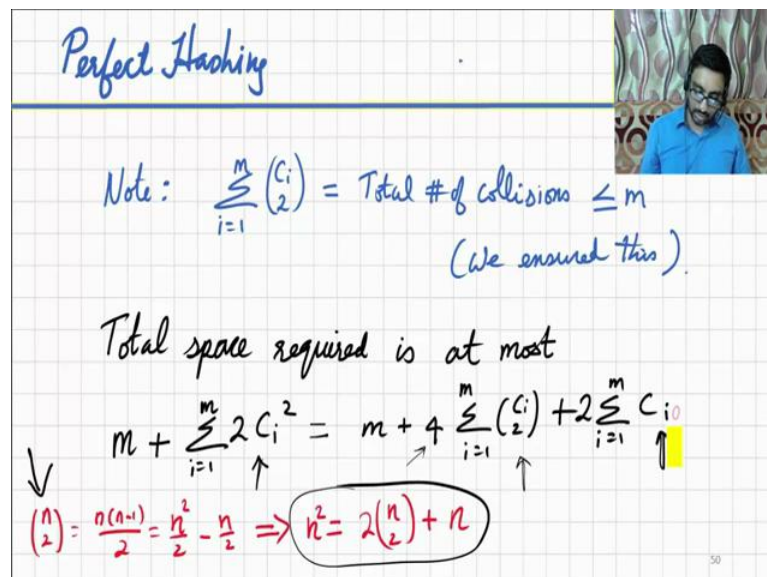
So, here is a picture. So, what we have is a hash table. Also, we use a hash function in h such that the number of collisions is kept to with m . We know that is possible and it can be done using the probabilistic method. We basically repeat until we find such a hash function. So, total number of collisions is at most m square. Summation over all i c_i square is at most m .

(Refer Slide Time: 11:00)



Let us focus on one particular hash location i at which some c_i items have been hashed into. So, change there is c_i items. What we do is, we take the chain of items and hash it into another hash table of length $2c_i^2$ and we ensure that this second hash table has no collisions what so ever.

(Refer Slide Time: 11:35)



As I just mentioned a little while ago if you count a total number of collisions basically summation i equal to 1 to m C_i choose 2 that is basically total number of collisions and that at most m . We ensure this. Now, Let us try to notice first of all that the search time is

$O(1)$. Let us say we have an item you want to search, we hash the item and look at its location within the hash table and if it hashes into a location with just one item in it or no items in it, then we are done.

We know that either the item is present or not present, but if we hash into a location that has collisions, so there is a chain at that location and then instead of going down the chain, we simply use in that location we insert. In fact, we do not even store the chain. We only need to store the function used to hash whether use to hash the items on the chain into another table, right. So, we take that hash function and again hash the same item and that will give us a location inside the second table and that could be empty or contain the item that we are looking for.

So, the query time is clear is of one question is what about the space required. Now, the complication comes from the fact that in addition to the original hash table at each location where multiple items that where hashed into, we took those items and we created a separate table for each one of them and now many such tables could be created.

So, the question is to ensure that the total space that we have used to create all of those tables is at most $O(m)$. So, Let us look at that. So, the total space required is that at most $m + \sum_{i=1}^m c_i^2$ because each time there was a collision at those items where there were items that where hashed into, we created separate hash tables of $2c_i^2$. So, that is the total space and that is equal to m plus this $2 \sum_{i=1}^m c_i^2$. We want to bring it to a form that we can analyze easily.



Recall that $\binom{n}{2}$ can be written as $\frac{n(n-1)}{2}$ of course, and that is equal to $\frac{n^2}{2} - \frac{n}{2}$. So, this means that n^2 can be written as two times $\binom{n}{2}$ plus n which is what we are going to do here. So, this $\sum_{i=1}^m c_i^2$ will going to write it as $2 \sum_{i=1}^m \binom{c_i}{2}$ and that is why we have 4 here plus two times the summation $\sum_{i=1}^m c_i$.

(Refer Slide Time: 15:12)

Perfect Hashing

Note: $\sum_{i=1}^m \binom{c_i}{2} = \text{Total \# of collisions} \leq m$
(We ensured this)

Total space required is at most

$$m + \sum_{i=1}^m 2c_i^2 = m + 4 \sum_{i=1}^m \binom{c_i}{2} + 2 \sum_{i=1}^m c_i$$
$$\leq m + 4m + 2m \in O(m)$$


51

Recall that we know that the summation over all i $\binom{c_i}{2}$ is at most m . So, this means that we can bound this entire quantity over here by $m + 4m + 2m$ which is at most $O(m)$.

With that we conclude our segment on perfect hashing. So, we have seen so far is that we can store a static dictionary using a hash table in such a manner that we can search for each item in $O(1)$ time and more importantly we do not need the simple uniform hashing assumption. For this instead all we need is that the hash functions we use are drawn uniformly at random from a family of hash functions from a two universal family of hash functions.