Data Structures and Algorithms (CS210A)

Lecture 40

- Search data structure for integers : Hashing
- Quick sort : some facts

Data structures for searching

in O(1) time

Motivating Example

Input: a given set *S* of **1009** positive integers

Aim: Data structure for searching

Example

```
{
123, 579236, 1072664, 770832456778, 61784523, 100004503210023, 19,
762354723763099, 579, 72664, 977083245677001238, 84, 100004503210023,
...
}
Data structure : Array storing S in sorted order
```

Searching : Binary search

O(log |S|) time



Problem Description



A search query:

Aim: A data structure for a <u>given</u> set *S* that can facilitate search in O(1) time

A trivial data structure for O(1) search time

Build a 0-1 array **A** of size **m** such that

A[i] = 1 if $i \in S$.

A[i] = 0 if $i \notin S$.

Time complexity for searching an element in set S: O(1).



This is a totally Impractical data structure because $n \ll m$! Example: n = few thousands, m = few trillions.

Question:

Can we have a data structure of O(n) size that can answer a search query in O(1) time?

Answer: Hashing

Hash function



Hash function:

h is a mapping from *U* to $\{0,1,...,n-1\}$ with the following characteristics.

- Space required for *h*
- *h(i)* computable in

Example:

Hash value:

For a given hash function h, and $i \in U$. h(i) is called hash value of i

Hash function, hash value



Hash function:

h is a mapping from **U** to $\{0, 1, ..., n - 1\}$ with the following characteristics.

- **Space** required for *h* : a few **words**.
- *h*(*i*) computable in O(1) time in word RAM.

Example: $h(i) = i \mod n$

Hash value:

For a given hash function h, and $i \in U$. h(i) is called hash value of i

Hash Table:

An array $T[0 \dots n - 1]$



Hash function:

h is a mapping from **U** to $\{0, 1, ..., n - 1\}$ with the following characteristics.

- **Space** required for *h* : a few **words**.
- **h**(*i*) computable in **O(1) time** in **word RAM.**

Example: $h(i) = i \mod n$

Hash value:

For a given hash function h, and $i \in U$. h(i) is called hash value of i

Hash Table:

An array $T[0 \dots n - 1]$



Hash function:

h is a mapping from **U** to $\{0, 1, ..., n - 1\}$ with the following characteristics.

- **Space** required for *h* : a few **words**.
- h(i) computable in **O(1) time** in word RAM.

Example: $h(i) = i \mod n$

Hash value:

For a given hash function h, and $i \in U$. h(i) is called hash value of i

Hash Table:

An array $T[0 \dots n - 1]$ of pointers storing **S**.



Question:

How to use (h,T) for searching an element $i \in U$? **Answer:** $k \leftarrow h(i)$; Search element i in the list T[k].

Time complexity for searching: O(length of the **longest** list in **T**).

Efficiency of Hashing depends upon hash function

A hash function *h* is <u>good</u> if it can **evenly** distributes *S*.

Aim: To search for a good hash function for a given set *S*.



There **can not be** any hash function **h** which is good for **every** *S***.**



For every h, there exists a subset of $\left[\frac{m}{n}\right]$ elements from U which are hashed to same value under h. So we can always construct a subset S for which all elements have same hash value

- \rightarrow All elements of this set S are present in a single list of the hah table T associated with h.
- \rightarrow O(n) worst case search time.

 $h(i) = i \bmod n$







 $h(x) = x \mod n$

m (

 $\left[\frac{m}{n}\right] - 1$

Let y_1, y_2, \dots, y_n denote n elements selected <u>randomly uniformly</u> from **U** to form **S**. **Question:** What is expected number of elements of S -i-ncolliding with y_1 ? **Answer:** Let y_1 takes value *i*. ●← <u>i</u> $P(y_i \text{ collides with } y_1) =$ -i+n $\left[\frac{m}{n}\right] - 1$ m-1i + 2nExpected number of elements of **S** colliding with $y_1 =$ -i + 3n $=\frac{\left[\frac{m}{n}\right]-1}{m-1}\left(n-1\right)$

$$= 0(1)$$

Conclusion

h(i) = i mod n works so well because
 for a uniformly random subset of U,
 the expected number of collision at an index of T is O(1).

It is easy to fool this hash function such that it achieves O(s) search time. (do it as a simple exercise).

This makes us think:

"How can we achieve worst case O(1) search time for a given set S."

Hashing: theory

 $U: \{0, 1, \dots, m-1\}$ $S \subseteq U,$ n = |S|,

Theorem [FKS, 1984]:

A hash table and hash function can be computed **Space**: O(n)**Query time**: worst case O(1)

Ingredients:

- elementary knowledge of prime numbers.
- The algorithms use **simple randomization**.

(We shall discuss such an algorithm in CS345.)





Quick Sort

Facts

(invented by Tony Hoare in 1960)

Quick sort versus Merge Sort

	Merge Sort	Quick Sort
Average case comparisons	$n \log_2 n$	1.39 $n \log_2 n$
Worst case comparisons	n log ₂ n	n(n-1)

Realization from Programming assignment 4 (part 1):

	<i>n</i> = 100	<i>n</i> = 1000	<i>n</i> ≥ 10000
No. of times Merge sort outperformed Quick sort	0 . 1 %	0 . 02 %	0%

Reasons:

- Overhead of **Copying** in merging **?**
- Technical (cache)

Very few students tried to find out 🔅

What makes Quick sort popular ?

No. of repetitions = **1000**

No. of times run time exceeds average by	100	1000	10 ⁴	10 ⁵	10 ⁶
10%	190	49	22	10	3
20%	28	17	12	3	0
50 %	2	1	1	0	0
100 %	0	0	0	0	0

Inference:

The chances of deviation from average case

→ The *reliability* of quick sort



What makes Quick sort popular ?



But a serious problem with Quick sort.

- Distribution sensitive 😕
- Can be fooled easily
 - sort in increasing order
 - Sort in decreasing order

Solution:

Select pivot element randomly uniformly in each call

This is randomized quick sort.