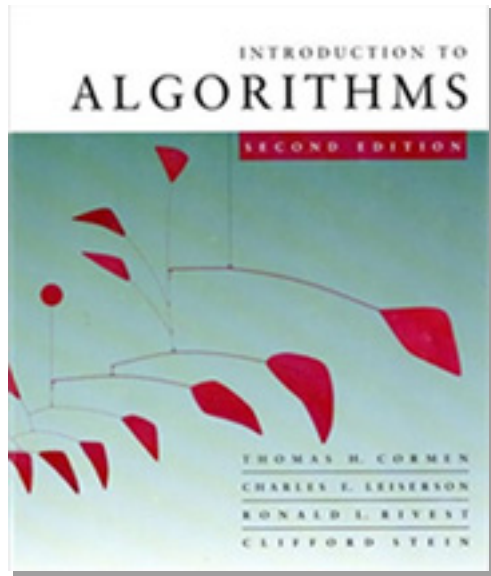


2.

Introduction to Algorithms

6.046J/18.401J



LECTURE 7

Hashing I

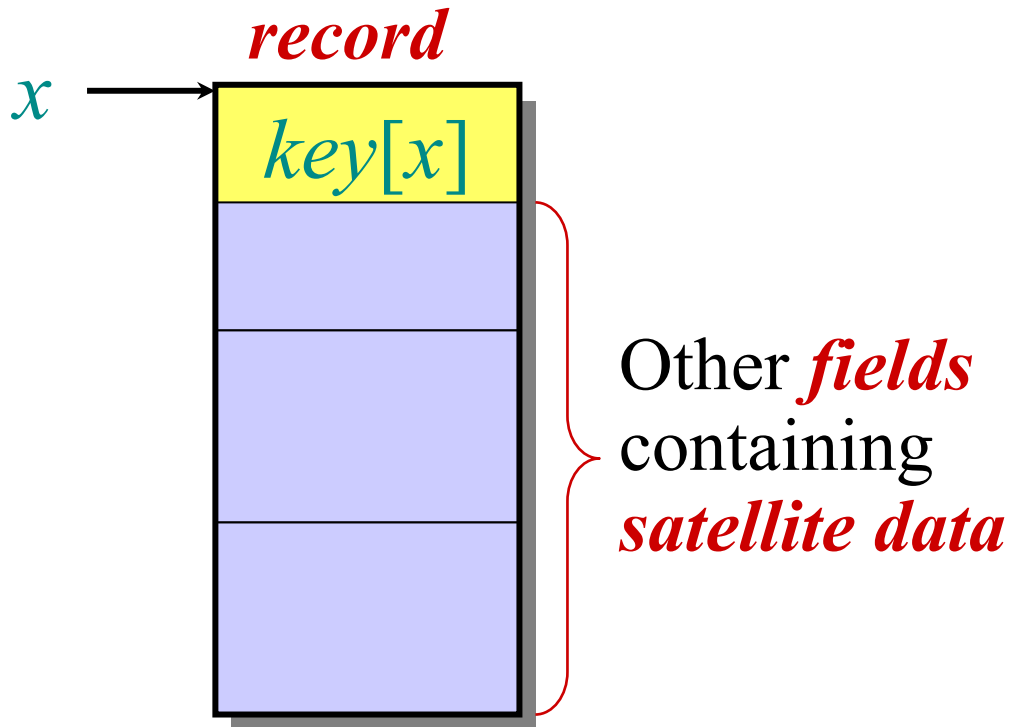
- Direct-access tables
- Resolving collisions by chaining
- Choosing hash functions
- Open addressing

Prof. Charles E. Leiserson



Symbol-table problem

Symbol table S holding n *records*:



Operations on S :

- $\text{INSERT}(S, x)$
- $\text{DELETE}(S, x)$
- $\text{SEARCH}(S, k)$

How should the data structure S be organized?



Direct-access table

IDEA: Suppose that the keys are drawn from the set $U \subseteq \{0, 1, \dots, m-1\}$, and keys are distinct. Set up an array $T[0 \dots m-1]$:

$$T[k] = \begin{cases} x & \text{if } x \in K \text{ and } \text{key}[x] = k, \\ \text{NIL} & \text{otherwise.} \end{cases}$$

Then, operations take $\Theta(1)$ time.

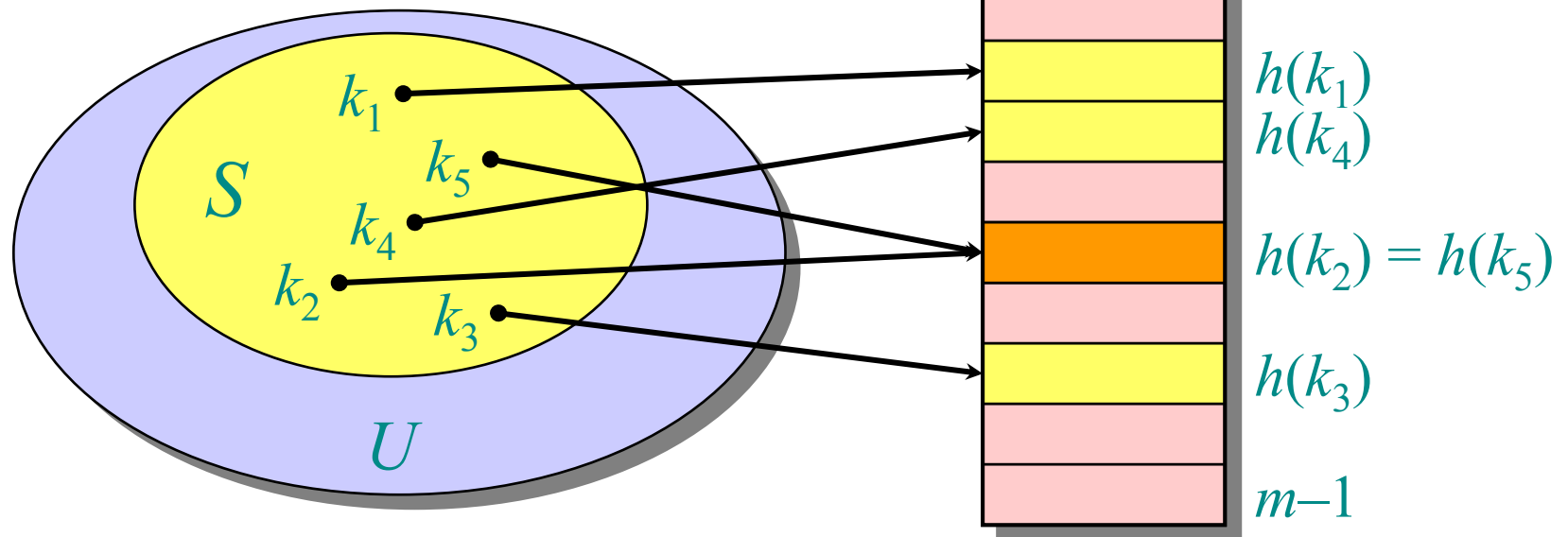
Problem: The range of keys can be large:

- 64-bit numbers (which represent 18,446,744,073,709,551,616 different keys),
- character strings (even larger!).

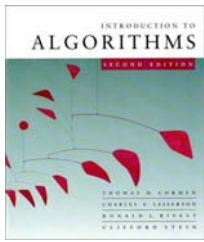


Hash functions

Solution: Use a *hash function* h to map the universe U of all keys into $\{0, 1, \dots, m-1\}$:

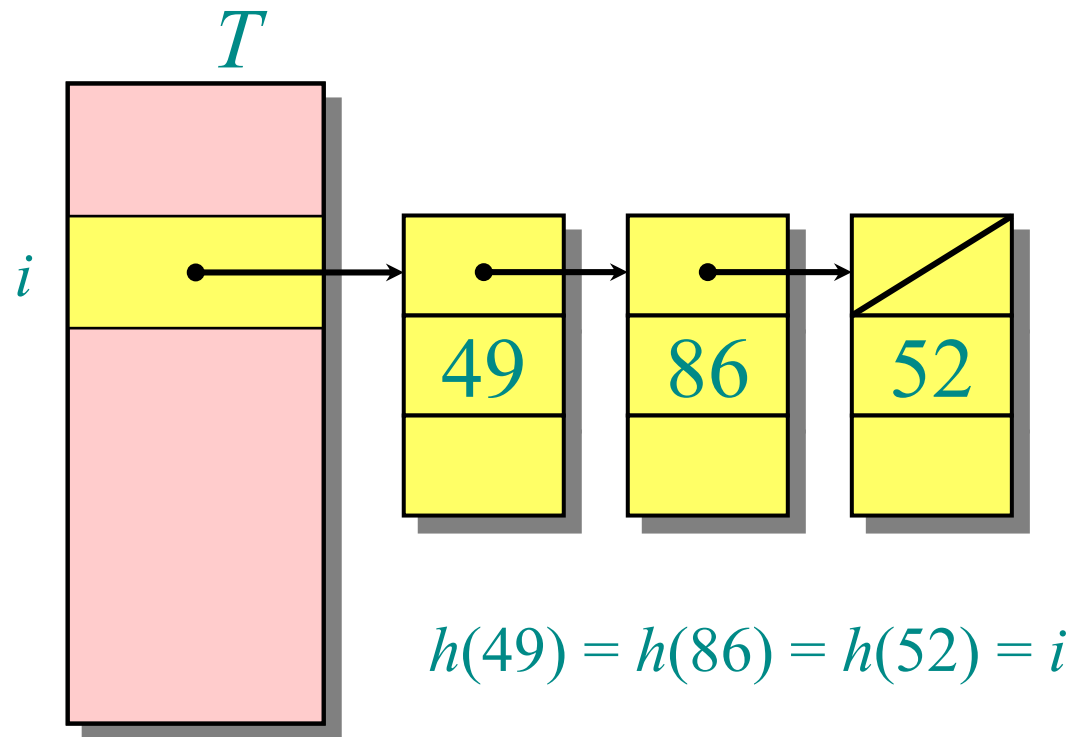


When a record to be inserted maps to an already occupied slot in T , a *collision* occurs.



Resolving collisions by chaining

- Link records in the same slot into a list.



Worst case:

- Every key hashes to the same slot.
- Access time = $\Theta(n)$ if $|S| = n$



Average-case analysis of chaining

We make the assumption of *simple uniform hashing*:

- Each key $k \in S$ is equally likely to be hashed to any slot of table T , independent of where other keys are hashed.

Let n be the number of keys in the table, and let m be the number of slots.

Define the *load factor* of T to be

$$\alpha = n/m$$

= average number of keys per slot.



Search cost

The expected time for an *unsuccessful* search for a record with a given key is $= \Theta(1 + \alpha)$.



Search cost

The expected time for an *unsuccessful* search for a record with a given key is

$$= \Theta(1 + \alpha).$$

*search
the list*

*apply hash function
and access slot*



Search cost

The expected time for an *unsuccessful* search for a record with a given key is

$= \Theta(1 + \alpha)$.

search the list (pointing to α)

apply hash function and access slot (pointing to 1)

Expected search time $= \Theta(1)$ if $\alpha = O(1)$, or equivalently, if $n = O(m)$.



Search cost

The expected time for an *unsuccessful* search for a record with a given key is

$= \Theta(1 + \alpha)$.

search the list (arrow pointing to α)

apply hash function and access slot (arrow pointing to 1)

Expected search time $= \Theta(1)$ if $\alpha = O(1)$, or equivalently, if $n = O(m)$.

A *successful* search has same asymptotic bound, but a rigorous argument is a little more complicated. (See textbook.)



Choosing a hash function

The assumption of simple uniform hashing is hard to guarantee, but several common techniques tend to work well in practice as long as their deficiencies can be avoided.

Desirata:

- A good hash function should distribute the keys uniformly into the slots of the table.
- Regularity in the key distribution should not affect this uniformity.



Division method

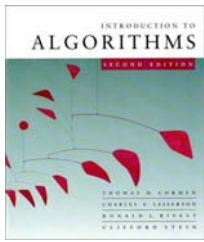
Assume all keys are integers, and define

$$h(k) = k \bmod m.$$

Deficiency: Don't pick an m that has a small divisor d . A preponderance of keys that are congruent modulo d can adversely affect uniformity.

Extreme deficiency: If $m = 2^r$, then the hash doesn't even depend on all the bits of k :

- If $k = 1011000111\underbrace{011010}_2$ and $r = 6$, then
$$h(k) = 011010_2.$$



Division method (continued)

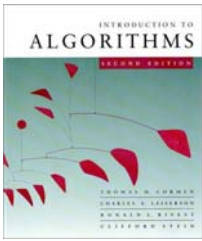
$$h(k) = k \bmod m.$$

Pick m to be a prime not too close to a power of 2 or 10 and not otherwise used prominently in the computing environment.

Annoyance:

- Sometimes, making the table size a prime is inconvenient.

But, this method is popular, although the next method we'll see is usually superior.



Multiplication method

Assume that all keys are integers, $m = 2^r$, and our computer has w -bit words. Define

$$h(k) = (A \cdot k \bmod 2^w) \text{ rsh } (w - r),$$

where **rsh** is the “bitwise right-shift” operator and A is an odd integer in the range $2^{w-1} < A < 2^w$.

- Don't pick A too close to 2^{w-1} or 2^w .
- Multiplication modulo 2^w is fast compared to division.
- The **rsh** operator is fast.