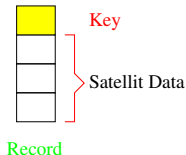# Hashing

AiC FME, UPC

Fall 2021

# Data Structures: Reminder

Given a universe $\mathcal{U}$, a dynamic set of records, where each record:



- Array
- Linked List (and variations)
- Stack (LIFO): Supports push and pop
- Queue (FIFO): Supports enqueue and dequeue
- Deque: Supports push, pop, enqueue and dequeue
- Heaps: Supports insertions, deletions, find Max and MIN
- Hashing

# Data structures for dynamic sets

**DICTIONARY**

Data structure for maintaining $\mathcal{S} \subset \mathcal{U}$ together with operations:

- Search($k$): decide if $k \in \mathcal{S}$
- Insert($k$): $\mathcal{S} := \mathcal{S} \cup \{k\}$
- Delete($k$): $\mathcal{S} := \mathcal{S} \backslash \{k\}$

**PRIORITY QUEUE**

Data structure for maintaining $\mathcal{S} \subset \mathcal{U}$ together with operations:

- Insert($x, k$): $\mathcal{S} := \mathcal{S} \cup \{x\}$
- Maximum(): Returns element of $\mathcal{S}$ with largest key value
- Extract-Maximum(): Returns $(x, k)$ with $k$ largest value in $\mathcal{S}$, $\mathcal{S} = \mathcal{S} - \{x\}$.

# Priority Queue implementations

**Linked List:**
- *INSERT:* $O(n)$
- *EXTRACT-MAX:* $O(1)$

**Heap:**
- *INSERT:* $O(\lg n)$
- *EXTRACT-MAX:* $O(\lg n)$

Using a Heap is a good compromise between fast insertion and slow extraction.

# Hashing

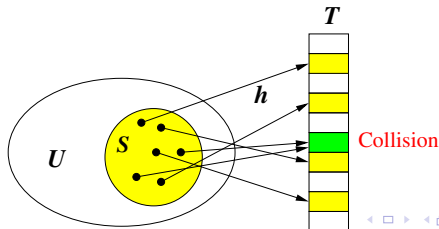Data Structure that supports *dictionary* operations on an universe of numerical keys.

Notice the number of possible keys represented as 64-bit integers is $2^{63} = 18446744073709551616$.

Tradeoff *time/space*

Define a hashing table $T[0, \ldots, m-1]$

a hashing function $h : \mathcal{U} \to T[0, \ldots, m-1]$

Hans P. Luhn
(1896-1964)

# Simple uniform hashing function.

- We want to store a maximum of $n$ keys in a hashing table $T$ with $m$ slots.
- The performance of hashing depends on how well $h$ distributes the keys on the $m$ slots.
- $h$ is simple uniform if it hash any key *with equal probability* into any slot, independently of where other keys go.
- In this way, we get a load factor $\alpha = n/m$, the average number of keys per slot.

# How to choose $h$?

Advice: For an exhaustive treaty on Hashing: D. Knuth, Vol. 3 of *The Art of computing programming*

# How to choose $h$?

Advice: For an exhaustive treaty on Hashing: D. Knuth, Vol. 3 of *The Art of computing programming*



$h$ depends on the type of key:

- For keys in the real interval $[0, 1)$, we can use $h(k) = \lfloor mk \rfloor$.

- For keys in the real interval $[s, t)$ scale by $1/(t - s)$, and use the previous method, $h(k/(t - s)) = \lfloor mk/(t - s) \rfloor$.

# The division method

Choose $m$ prime or as far as possible from a power of 2,

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

Fast ($\Theta(1)$) to compute in most languages ($k \% m$)!

Be aware: if $m = 2^r$ the hash does not depend on all the bits of K

If $r = 6$ with $k = 1011000111\underbrace{011010}_{=h(k)}$

($45530 \mod 64 = 858 \mod 64$)

In some applications, the keys may be very large, for instance with alphanumeric keys, which must be converted to ascii, and reinterpreted as numbers in binary.

Example: *averylongkey* is converted via ascii:

$97 \cdot 128^{11} + 118 \cdot 128^{10} +$
$101 \cdot 128^{9} + 114 \cdot 128^{8}$
$+121 \cdot 128^{7} + 108 \cdot 126^{6}$
$+111 \cdot 128^{5} + 110 \cdot 128^{4}$
$+103 \cdot 128^{3} + 107 \cdot 128^{2}$
$+101 \cdot 128^{1} + 121 \cdot 128^{0} = n$

which has 84-bits!



Source: www.LookupTables.com

# How to deal with large $n$?

For large $n$, to compute $h = n \mod m$, we can use mod arithmetic $+$ Horner's method:

$$(((((((((97 \cdot 128 + 118) \cdot 128 + 101) \cdot 128 + 114) \cdot 128 + 121)$$
$$\cdot 128 + 111) \cdot 128 + 110) \cdot 128 + 103) \cdot 128 + 107)$$
$$\cdot 128 + 101) \cdot 128 + 121 \quad \mod m$$
$$= (((((((((\underbrace{97 \cdot 128 + 118 \quad \mod m) \cdot 128) \quad \mod m}_{} + 101) \cdot \ldots) ))))))$$

# Collision resolution: Separate chaining

For each table address, construct a linked list of the items whose keys hash to that address.

- Every key goes to the same slot
- Time to explore the list = length of the list



$h(20)=h(27)=h(8)=i$

# Cost of average analysis of chaining

The cost of the dictionary operations using hashing:

- Insertion of a new key: $\Theta(1)$.
- Search of a key: $O($ length of the list$)$
- Deletion of a key: $O($length of the list$)$.

Under the hypothesis that $h$ is *simply uniform hashing*, each key $x$ is equally likely to be hashed to any slot of $T$, independently of where other keys are hashed

Therefore, the expected number of keys falling into $T[i]$ is $\alpha = n/m$.

# Cost of search

- For an unsuccessful search ($x$ is not in $T$), we have to explore the ist at $h(x) \to T[i]$. So, the expected time to search the list at $T[i]$ is $O(1 + \alpha)$.
  ($\alpha$ of searching the list and $\Theta(1)$ of computing $h(x)$ and going to slot $T[i]$)

# Cost of search

- For an unsuccessful search ($x$ is not in $T$), we have to explore the ist at $h(x) \to T[i]$. So, the expected time to search the list at $T[i]$ is $O(1 + \alpha)$.
  ($\alpha$ of searching the list and $\Theta(1)$ of computing $h(x)$ and going to slot $T[i]$)

- For an successful search, we obtain the same bound, although in most of the cases we would have to search a fraction of the list until finding the $x$ element.)

# Cost of search

- For an unsuccessful search ($x$ is not in $T$), we have to explore the ist at $h(x) \to T[i]$. So, the expected time to search the list at $T[i]$ is $O(1 + \alpha)$.
  ($\alpha$ of searching the list and $\Theta(1)$ of computing $h(x)$ and going to slot $T[i]$)

- For an successful search, we obtain the same bound, although in most of the cases we would have to search a fraction of the list until finding the $x$ element.)

- Under the assumption of simple uniform hashing, in a hash table with chaining, a search takes time $\Theta(1 + \frac{n}{m})$ on average.

# Cost of search

- For an unsuccessful search ($x$ is not in $T$), we have to explore the ist at $h(x) \rightarrow T[i]$. So, the expected time to search the list at $T[i]$ is $O(1 + \alpha)$.
  ($\alpha$ of searching the list and $\Theta(1)$ of computing $h(x)$ and going to slot $T[i]$)

- For an successful search, we obtain the same bound, although in most of the cases we would have to search a fraction of the list until finding the $x$ element.)

- Under the assumption of simple uniform hashing, in a hash table with chaining, a search takes time $\Theta(1 + \frac{n}{m})$ on average.

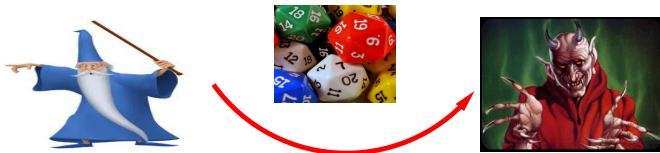- Notice that if $n = \theta(m)$ then $\alpha = O(1)$ and search time is $\Theta(1)$.

# Universal hashing: Motivation



- For every deterministic hash function, there is a set of bad instances.
- An adversary can arrange the keys so your function hashes most of them to the same slot.

# Universal hashing: Motivation



- For every deterministic hash function, there is a set of bad instances.
- An adversary can arrange the keys so your function hashes most of them to the same slot.
- Create a set $\mathcal{H}$ of hash functions on $\mathcal{U}$ and choose a hashing function at random and independently of the keys.
- The adversary might known the probability space but not the particular selection.

# Universal hashing

Let $\mathcal{U}$ be the universe of keys and let $\mathcal{H}$ be a collection of hashing functions with hashing table $T[0, \ldots, m-1]$, $\mathcal{H}$ is universal if $\forall x, y \in \mathcal{U}, x \neq y$, then

$$|\{h \in \mathcal{H} \mid h(x) = h(y)\}| \leq \frac{|\mathcal{H}|}{m}.$$

In an equivalent way, $\mathcal{H}$ is *universal* if $\forall x, y \in \mathcal{U}, x \neq y$, and for any $h$ chosen uniformly from $\mathcal{H}$, we have

$$\mathbf{Pr}\left[h(x) = h(y)\right] \leq \frac{1}{m}.$$

# Universality gives good average-case behaviour

Theorem

*If we pick u.a.r. h from a universal family $\mathcal{H}$ and build a table with size m for a set of n keys, for any given key x let $C_x$ be a random variable counting the number of collisions with others keys y in T.*

$$\mathbf{E}\left[C_x\right] \le n/m.$$

# Construction of a universal family: $\mathcal{H}$

Let $\mathcal{U}$ be the key universe and let $N$ be the maximum key value. Our target is a hash table with $m$ positions, $T[0, \ldots, m-1]$.

- *Choose a prime $p$, $N \leq p \leq 2N$. Then $\mathcal{U} \subset \mathbb{Z}_p = \{0, 1, \ldots, p-1\}$.*
- Define $\mathcal{H} = \{h_{a,b} | a, b \in \mathbb{Z}_p, a \neq 0\}$.

# Construction of a universal family: $\mathcal{H}$

Let $\mathcal{U}$ be the key universe and let $N$ be the maximum key value. Our target is a hash table with $m$ positions, $T[0, \ldots, m-1]$.

- *Choose a prime $p$, $N \leq p \leq 2N$. Then $\mathcal{U} \subset \mathbb{Z}_p = \{0, 1, \ldots, p-1\}$.*
- Define $\mathcal{H} = \{h_{a,b} | a, b \in \mathbb{Z}_p, a \neq 0\}$.

- To select u.a.r. $h \in \mathcal{H}$, *choose independently and u.a.r. $a \in \mathbb{Z}_p^+$ and $b \in \mathbb{Z}_p$. Given a key $x$ define $h_{a,b}(x) = (\underbrace{(ax + b) \mod p}_{g_{a,b}(x)}) \mod m$.*

# Construction of a universal family: $\mathcal{H}$

Let $\mathcal{U}$ be the key universe and let $N$ be the maximum key value. Our target is a hash table with $m$ positions, $T[0, \ldots, m-1]$.

- *Choose a prime $p$, $N \leq p \leq 2N$. Then $\mathcal{U} \subset \mathbb{Z}_p = \{0, 1, \ldots, p-1\}$.*
- Define $\mathcal{H} = \{h_{a,b} | a, b \in \mathbb{Z}_p, a \neq 0\}$.

- To select u.a.r. $h \in \mathcal{H}$, *choose independently and u.a.r. $a \in \mathbb{Z}_p^+$ and $b \in \mathbb{Z}_p$. Given a key $x$ define $h_{a,b}(x) = (\underbrace{(ax + b) \mod p}_{g_{a,b}(x)}) \mod m$.*

- Example: $p = 17, m = 6$, we have $\mathcal{H}_{17,6} = \{h_{a,b} : a \in \mathbb{Z}_p^+, b \in \mathbb{Z}_p\}$
  if $x = 8, a = 3, b = 4$ then
  $h_{3,4}(8) = ((3 \cdot 8 + 4) \mod 17) \mod 6 = 5$

# Properties of $\mathcal{H}$

1. $h_{ab} : \mathbb{Z}_p \to \mathbb{Z}_m$.
2. $|\mathcal{H}| = p(p-1)$. (We can select $a$ in $p-1$ ways and $b$ in $p$ ways)
3. Specifying an $h \in \mathcal{H}$ requires $O(\lg p) = O(\lg N)$ bits.
4. To choose $h \in \mathcal{H}$ select $a, b$ independently and u.a.r. from $\mathbb{Z}_p^+$ and $\mathbb{Z}_p$.
5. Evaluating $h(x)$ is fast.

Theorem

*The family $\mathcal{H}$ is universal.*

For the proof:
Chapter 11 of Cormen. Leiserson, Rivest, Stein: *An introduction to Algorithms*