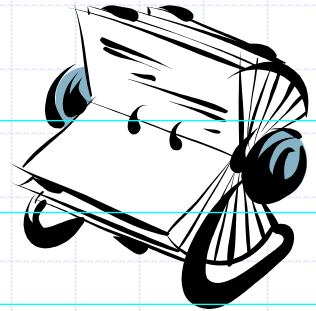


# Maps

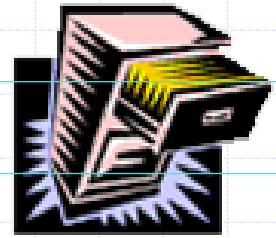


# Maps



- ◆ A map models a searchable collection of key-value entries
- ◆ The main operations of a map are for searching, inserting, and deleting items
- ◆ Multiple entries with the same key are **not** allowed [key = name value = address]
- ◆ Applications: →
  - address book
  - student-record database → [key = roll no, value = student data]

$$\{ k_1 \Rightarrow v_1, k_2 \Rightarrow v_2 \dots k_n \Rightarrow v_n \}$$



# The Map ADT (§ 8.1)

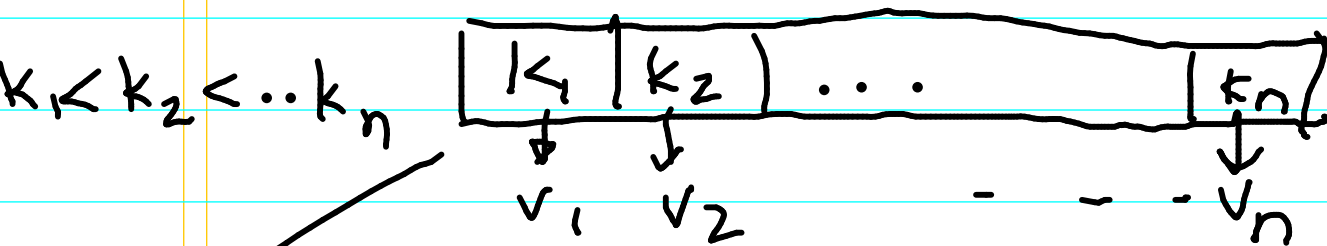
## ◆ Map ADT methods:

- **get**(k): if the map M has an entry with key k, return its associated value; else, return null
- **put**(k, v): insert entry (k, v) into the map M; if key k is not already in M, then return null; else, return old value associated with k  
*replace the value*
- **remove**(k): if the map M has an entry with key k, remove it from M and return its associated value; else, return null  
*No repetitions*
- **size()**, **isEmpty()**
- **keys()**: return an iterator of the keys in M  
*can have repetitions*
- **values()**: return an iterator of the values in M

$\{ 1500 \Rightarrow \text{"Bombay"} \}$  . put(1500, "Mumbai")  
will replace "Bombay" with "Mumbai" & return "Bombay"

Q: Implementation(s) of Map that take less time for `get(k)` while being "reasonable" for `put(k, v)`

Ideas: (a) Assuming a way of sorting keys



`get(k)` through binary search in  $O(\log(n))$

`put(k, v)` while respecting sorting will be  $O(n)$  (even with amortised analysis)

Eg: 

```
public class Pair {  
    String key;  
    String value;  
}
```

could implement as an array `Pair[]` but maintain sorting on key field

(b) Linked list of Pair:

$\text{get}(k) : O(n)$

$\text{put}(k, v) : O(1)$

$O(n)$  for a "faithful"

implementation that "gets"  $\langle k, v \rangle$  if  $k$  was already present

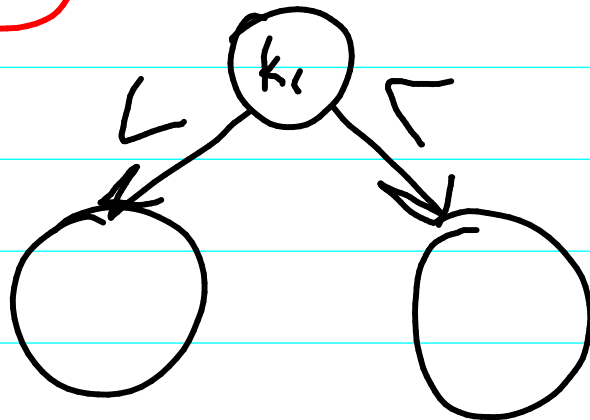
Assuming insertion at beginning so that  $\text{get}(k)$  gets last inserted value scanning from "front"

(c) Encoding function on key to map "key" to a "position" in an array

Q: Is this a "unique" position?

Can  $f(k_1) = f(k_2)$  for  $k_1 \neq k_2$  [i.e. collision for some  $k_1, k_2$ ]?

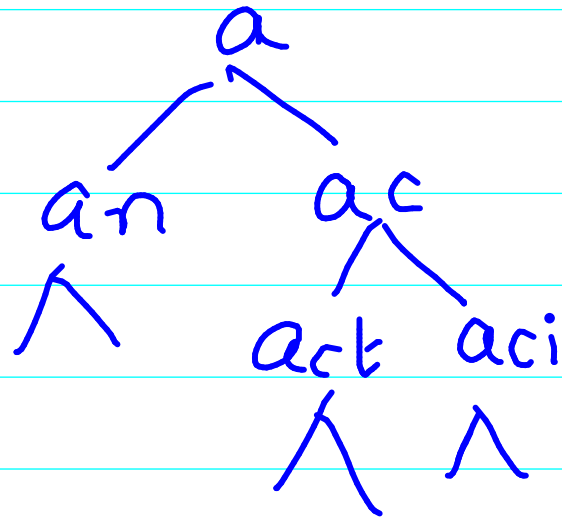
(d) Binary search tree



$\text{get}(k) : O(\log n)$

$\text{put}(k, v) : O(\log n)$

e) Trie (prefix tree)



f) Assuming key length  $\leq r$

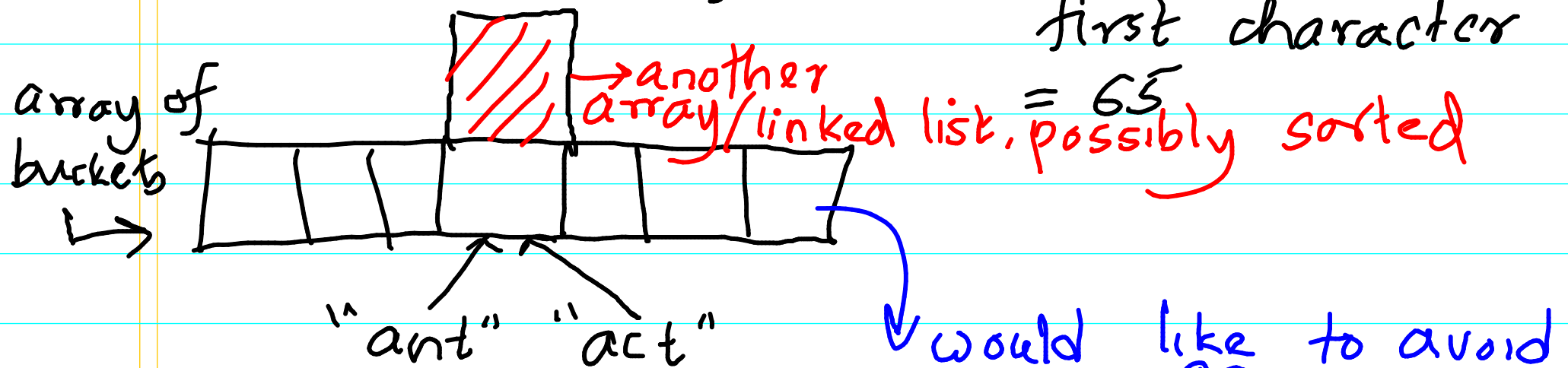
Take an " $r$ " dimensional array



## More elaboration on (c):

- ① Looks difficult to come up with a "hash" function  $f$  that will give no collisions.

Let us consider  $f(\text{"ant"}) = \text{ASCII value of first character}$



② Choice of "hash" functions

① Should give close to "uniform" distribution over all keys for the buckets

# Example

<i><b>Operation</b></i>	<i><b>Output</b></i>	<i><b>Map</b></i>
isEmpty()	<b>true</b>	$\emptyset$
put(5,A)	<b>null</b>	(5,A)
put(7,B)	<b>null</b>	(5,A),(7,B)
put(2,C)	<b>null</b>	(5,A),(7,B),(2,C)
put(8,D)	<b>null</b>	(5,A),(7,B),(2,C),(8,D)
put(2,E)	<i>C</i>	(5,A),(7,B),(2,E),(8,D)
get(7)	<i>B</i>	(5,A),(7,B),(2,E),(8,D)
get(4)	<b>null</b>	(5,A),(7,B),(2,E),(8,D)
get(2)	<i>E</i>	(5,A),(7,B),(2,E),(8,D)
size()	4	(5,A),(7,B),(2,E),(8,D)
remove(5)	<i>A</i>	(7,B),(2,E),(8,D)
remove(2)	<i>E</i>	(7,B),(8,D)
get(2)	<b>null</b>	(7,B),(8,D)
isEmpty()	<b>false</b>	(7,B),(8,D)



# Comparison to java.util.Map

## Map ADT Methods

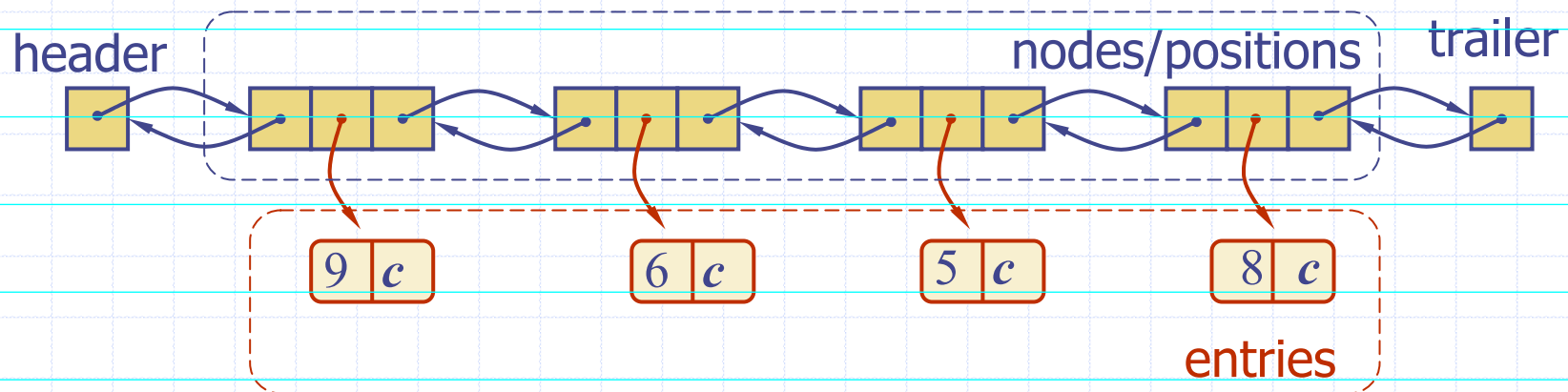
size()  
isEmpty()  
get( $k$ )  
put( $k, v$ )  
remove( $k$ )  
keys()  
values()

## java.util.Map Methods

size()  
isEmpty()  
get( $k$ )  
put( $k, v$ )  
remove( $k$ )  
keySet().iterator()  
values().iterator()

# A Simple List-Based Map

- ◆ We can efficiently implement a map using an unsorted list
  - We store the items of the map in a list  $S$  (based on a doubly-linked list), in arbitrary order



# The get(k) Algorithm

**Algorithm** get( $k$ ):

$B = S.positions()$  { $B$  is an iterator of the positions in  $S$ }

**while**  $B.hasNext()$  **do**

$p = B.next()$  the next position in  $B$

**if**  $p.element().key() = k$  **then**

**return**  $p.element().value()$

**return null** {there is no entry with key equal to  $k$ }

# The put( $k, v$ ) Algorithm

**Algorithm** put( $k, v$ ):

$B = S.positions()$

**while**  $B.hasNext()$  **do**

$p = B.next()$

**if**  $p.element().key() = k$  **then**

$t = p.element().value()$

$B.replace(p, (k, v))$

**return**  $t$                    {return the old value}

$S.insertLast((k, v))$

$n = n + 1$                    {increment variable storing number of entries}

**return null**               {there was no previous entry with key equal to  $k$ }

# The remove( $k$ ) Algorithm

**Algorithm** remove( $k$ ):

$B = S.positions()$

**while**  $B.hasNext()$  **do**

$p = B.next()$

**if**  $p.element().key() = k$  **then**

$t = p.element().value()$

$S.remove(p)$

$n = n - 1$

**return**  $t$

**return** null

{decrement number of entries}

{return the removed value}

{there is no entry with key equal to  $k$ }

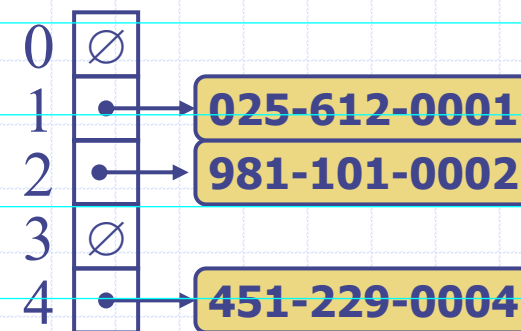
# Performance of a List-Based Map

## ◆ Performance:

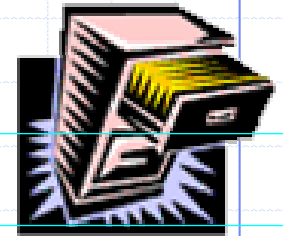
- **put** takes  $O(1)$  time since we can insert the new item at the beginning or at the end of the sequence
- **get** and **remove** take  $O(n)$  time since in the worst case (the item is not found) we traverse the entire sequence to look for an item with the given key

◆ The unsorted list implementation is effective only for maps of small size or for maps in which puts are the most common operations, while searches and removals are rarely performed (e.g., historical record of logins to a workstation)

# Hash Tables



# Recall the Map ADT (§ 8.1)



## ◆ Map ADT methods:

- **get(k)**: if the map  $M$  has an entry with key  $k$ , return its associated value; else, return null
- **put(k, v)**: insert entry  $(k, v)$  into the map  $M$ ; if key  $k$  is not already in  $M$ , then return null; else, return old value associated with  $k$
- **remove(k)**: if the map  $M$  has an entry with key  $k$ , remove it from  $M$  and return its associated value; else, return null
- **size()**, **isEmpty()**
- **keys()**: return an iterator of the keys in  $M$
- **values()**: return an iterator of the values in  $M$



U is universe of keys.  
S is evolving set of keys

[1] Array based:  $O(1)$   
operations but  $O(|U|)$  space  
[2] list based:  $O(|S|)$  space but  
 $O(|S|)$  lookup  
Desire:  $h: U \rightarrow \{0, 1, 2, \dots, n-1\}$

Q: What to do with collisions

Ignore leap year:

## THE BIRTHDAY PARADOX

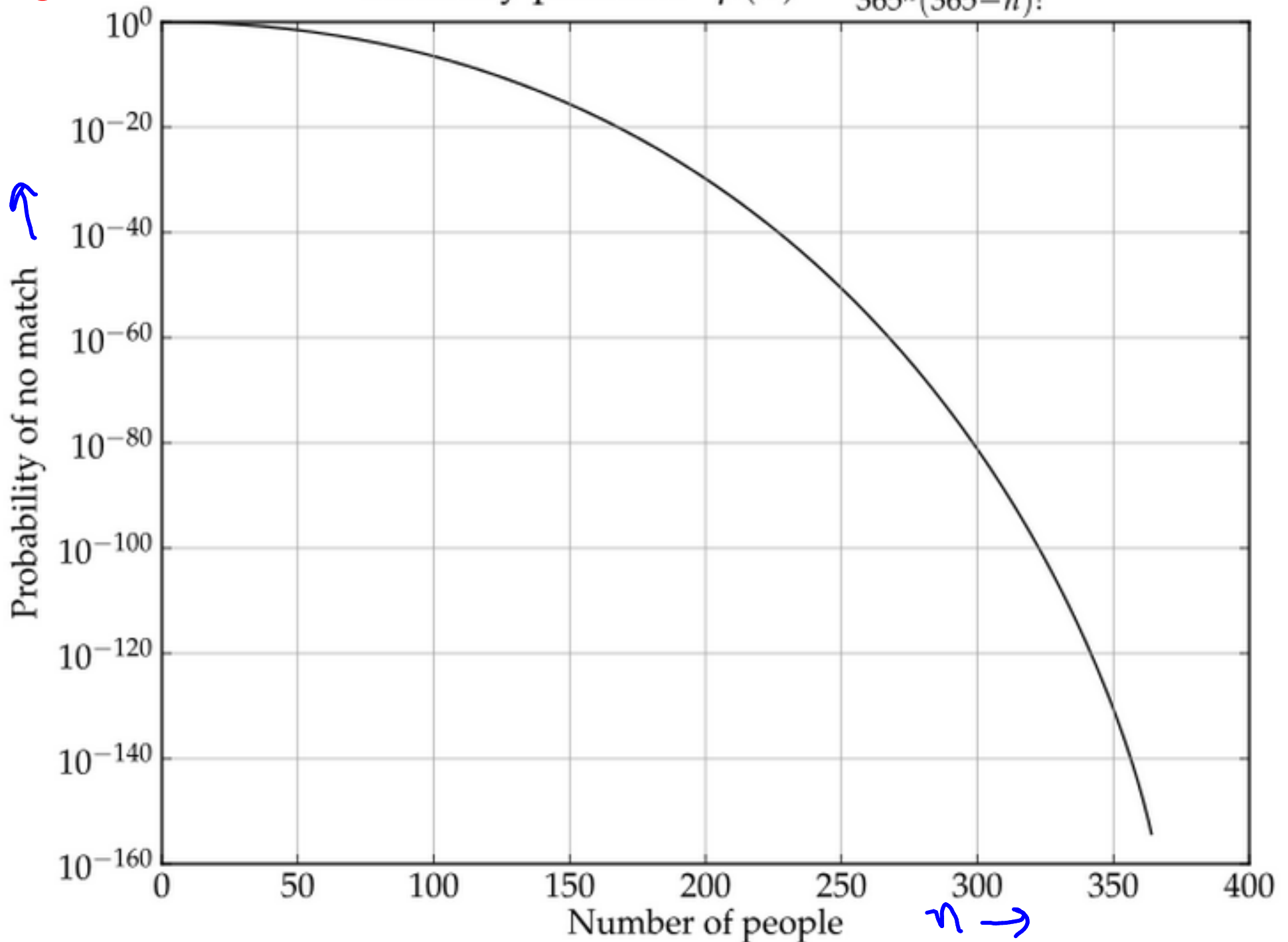
Consider  $n$  people with random birthdays (i.e. with each day of the year equally likely). How large does  $n$  need to be before there is at least a 50% chance that two people have the same birthday?

Choose: 23, 57, 184, 367

$$\begin{array}{ccccccc} & \underbrace{50\%} & \underbrace{99\%} & \underbrace{99.99\%} & \underbrace{100\%} & & \\ P(\text{amongst } n \text{ people at least 2 have same bday}) & & & & & & \\ = 1 - P(\text{all } n \text{ people have diff bdays}) & = & 1 - \frac{365 P_n}{(365)^n} & & & & \end{array}$$

<http://en.wikipedia.org/wiki/File:Birthdaymatch.svg>

Birthday paradox -  $\bar{p}(n) = \frac{365!}{365^n(365-n)!}$



# of people ( $n$ ) = # of possible keys

# of days in a year (365) = # of slots

that keys will be  
mapped to using  
hash function

Good hash fn



Best (non-practical)  
hash function (in terms  
of minimum collisions)  
is RANDOM

both  
reqd

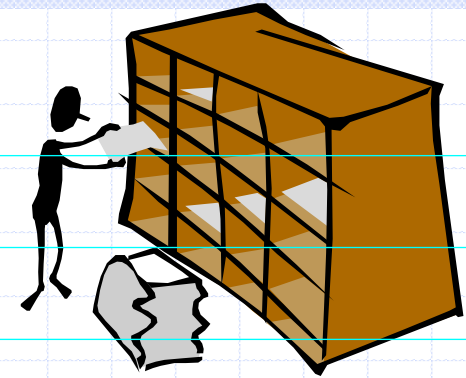
Deterministic

Strategy for  
collision handling



In practice we  
need deterministic  
hash fns. that  
closely mimic random

# Hash Functions and Hash Tables (§ 8.2)



- ◆ A **hash function**  $h$  maps keys of a given type to integers in a fixed interval  $[0, N-1]$

- ◆ Example:

$$h(x) = x \bmod N$$

is a hash function for integer keys

- ◆ The integer  $h(x)$  is called the **hash value** of key  $x$

- ◆ A **hash table** for a given key type consists of
  - Hash function  $h$
  - Array (called table) of size  $N$

- ◆ When implementing a map with a hash table, the goal is to store item  $(k, o)$  at index  $i = h(k)$

→ If  $x = 201$  &  $N = 10$   
 $h(x) = 1$

- ◆ A **hash table** for:  $U \rightarrow [0, N-1]$  consists of
  - Hash function  $h$
  - Array (called table) of size  $N$

→ Not  $O(N)$

(1) If  $x$  (keys) are distributed uniformly over all integers, then choice of  $N$  does not affect distribution. You will only choose  $N$  based on how much space/resources are available

(2) If  $x$  (keys) are not all equally possible (for eg: names of people have some patterns in a region) then  $x \bmod N$  could lead to unequal slot occupancy. Eg:  $x$  only evens &  $N=10$   
 $\therefore$  In such a case, the "uniform" hashing is ideal & realised, for eg by  $N = \text{large prime}$

$x = \text{all even}$

$$N=10 \Rightarrow x \bmod N \in \{0, 2, 4, 6, 8\}$$

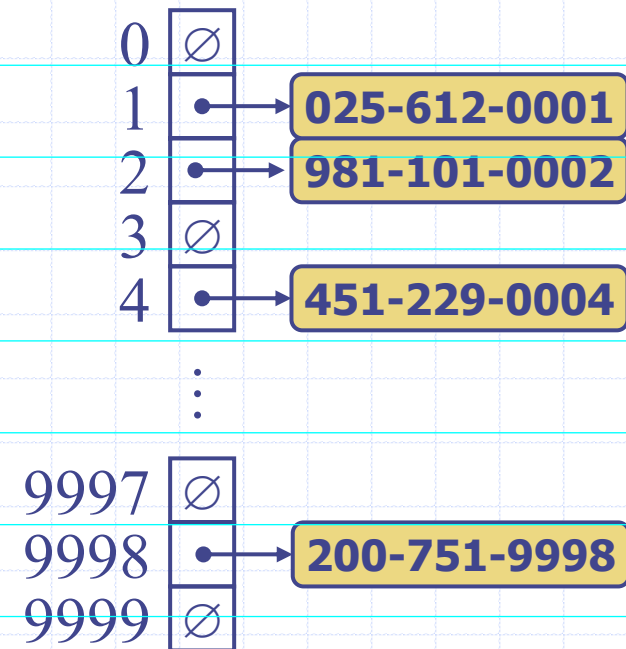
$$N=11 \Rightarrow x \bmod N \in \underbrace{\{0, 1, 2, \dots, 9\}}_{10}$$

$$N=2 \Rightarrow x \bmod N \in \{0, 1\}$$

Primes are good. To reduce collisions,  
large prime preferred

# Example

- ◆ We design a hash table for a map storing entries as (SSN, Name), where SSN (social security number) is a nine-digit positive integer
- ◆ Our hash table uses an array of size  $N = 10,000$  and the hash function  $h(x) = \text{last four digits of } x$





Application:

1] IP address lookup at a router requires hash table..

2] The 2-sum problem (repeated lookups) in constant time!

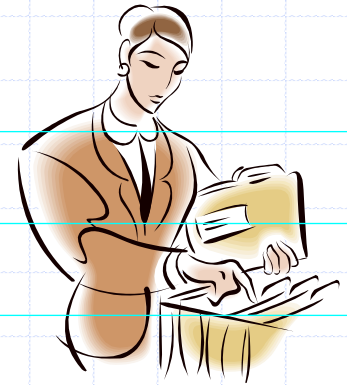
\* In general, for search algorithms: search tree algorithms, searching on graphs, chess playing program (search based configuration to find next possible moves)

} H/w  
Similar to histogram

3] Lookups for compilers! Symbol tables!

4] Blocking network traffic

# Hash Functions (§ 8.2.2)



- ◆ A hash function is usually specified as the composition of two functions.

*ex: Think of possibilities*

**Hash code:**

$h_1: \text{keys} \rightarrow \text{integers}$

**Compression function:**

$h_2: \text{integers} \rightarrow [0, N-1]$   
 $h_2 = x \bmod N$

- ◆ The hash code is applied first, and the compression function is applied next on the result, i.e.,

$$h(x) = h_2(h_1(x))$$

- ◆ The goal of the hash function is to “disperse” the keys in an apparently random way

*where N is a large prime is frequent choice*

*Neither needs to be one-one*

Example hash codes:

- ① Binary representation of the "key" object  
Eg: "serialised object" in java.
- ② The "binary" address of the object
- ③ For strings: Sum of ASCII values of characters.
- ④ For strings: Since "character sequences" are often characteristic in strings (eg: "qz" is unlikely to occur), position dependent encoding of characters

$$\text{Eg: } \sum_{i=1}^{|S|} \text{Ascii}(C_i) * N_1^{i-1}$$

$$\text{Where } S = C_1 C_2 \dots C_{|S|}$$

(Imagine writing the string with  $N_1 = 26$  in base 26)

$$\text{Eg: } N_1 = 26$$

$$h_1: S \xrightarrow{N_1} x$$

$$h_2: x \xrightarrow{x \bmod N_2} j \in [0, N_2]$$

Q: Will some relation between  $N_1$  &  $N_2$  help in practice? ANS:  $N_1$  &  $N_2$  coprime

Eg:  $N_2 = 8$  &  $x = \text{multiples of } 3 \Rightarrow$  all 8 slots filled up!

# Hash Codes (§ 8.2.3)



## ◆ Memory address:

- We reinterpret the memory address of the key object as an integer (default hash code of all Java objects)
- Good in general, except for numeric and string keys

## ◆ Integer cast:

- We reinterpret the bits of the key as an integer
- Suitable for keys of length less than or equal to the number of bits of the integer type (e.g., byte, short, int and float in Java)

## ◆ Component sum:

- We partition the bits of the key into components of fixed length (e.g., 16 or 32 bits) and we sum the components (ignoring overflows)
- Suitable for numeric keys of fixed length greater than or equal to the number of bits of the integer type (e.g., long and double in Java)

# Hash Codes (cont.)

## ◆ Polynomial accumulation:

- We partition the bits of the key into a sequence of components of fixed length (e.g., 8, 16 or 32 bits)

$$a_0 a_1 \dots a_{n-1}$$

- We evaluate the polynomial  $p(z) = a_0 + a_1 z + a_2 z^2 + \dots + a_{n-1} z^{n-1}$

at a fixed value  $z$ , ignoring overflows

- Especially suitable for strings (e.g., the choice  $z = 33$  gives at most 6 collisions on a set of 50,000 English words)

## ◆ Polynomial $p(z)$ can be evaluated in $O(n)$ time using Horner's rule:

- The following polynomials are successively computed, each from the previous one in  $O(1)$  time

$$p_0(z) = a_{n-1}$$

$$p_i(z) = a_{n-i-1} + z p_{i-1}(z) \\ (i = 1, 2, \dots, n-1)$$

## ◆ We have $p(z) = p_{n-1}(z)$

such as names/ids, char positions imp.

Experimental results: Imp for assignment 6

Optimise multiplications

# Compression Functions

## (§ 8.2.4)



### ◆ Division:

- $h_2(y) = y \bmod N$
- The size  $N$  of the hash table is usually chosen to be a prime
- The reason has to do with number theory and is beyond the scope of this course

If polynomial encoding (accumulation)  
then prefer:  $\mathbb{Z}$  &  $N$   
are co prime

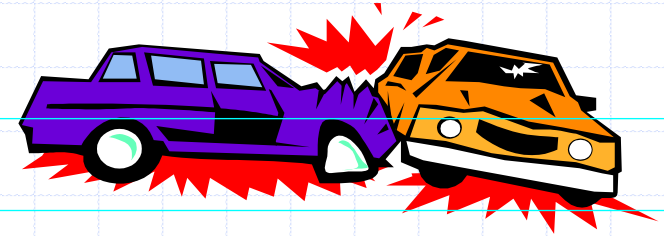
### ◆ Multiply, Add and Divide (MAD):

- $h_2(y) = (ay + b) \bmod N$
- $a$  and  $b$  are nonnegative integers such that  $a \bmod N \neq 0$
- Otherwise, every integer would map to the same value  $b$

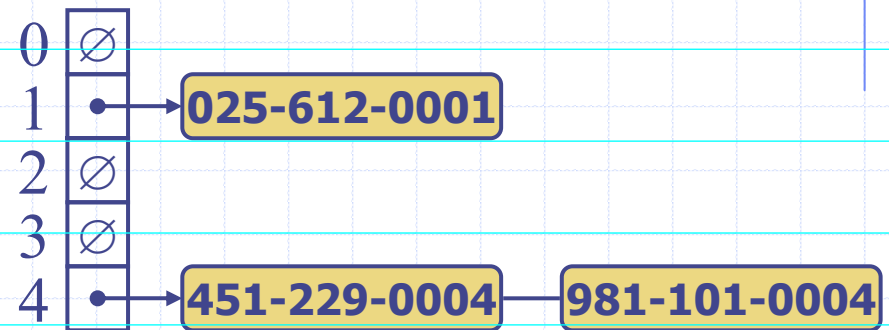
else  
will be  
meaningless

# Collision Handling

## (§ 8.2.5)



- ◆ Collisions occur when different elements are mapped to the same cell



- ◆ **Separate Chaining:**  
let each cell in the table point to a linked list of entries that map there

- ◆ Separate chaining is simple, but requires additional memory (pointers) outside the table



# Map Methods with Separate Chaining used for Collisions

◆ Delegate operations to a list-based map at each cell:

**Algorithm** get( $k$ ):

**Output:** The value associated with the key  $k$  in the map, or **null** if there is no entry with key equal to  $k$  in the map

**return**  $A[h(k)].get(k)$

*Involving one data structure within another*  
{delegate the get to the list-based map at  $A[h(k)]$ }

**Algorithm** put( $k, v$ ):

**Output:** If there is an existing entry in our map with key equal to  $k$ , then we return its value (replacing it with  $v$ ); otherwise, we return **null**

$t = A[h(k)].put(k, v)$

{delegate the put to the list-based map at  $A[h(k)]$ }

**if**  $t = \text{null}$  **then**

{ $k$  is a new key}

$n = n + 1$

**return**  $t$

**Algorithm** remove( $k$ ):

**Output:** The (removed) value associated with key  $k$  in the map, or **null** if there is no entry with key equal to  $k$  in the map

$t = A[h(k)].remove(k)$

{delegate the remove to the list-based map at  $A[h(k)]$ }

**if**  $t \neq \text{null}$  **then**

{ $k$  was found}

$n = n - 1$

**return**  $t$

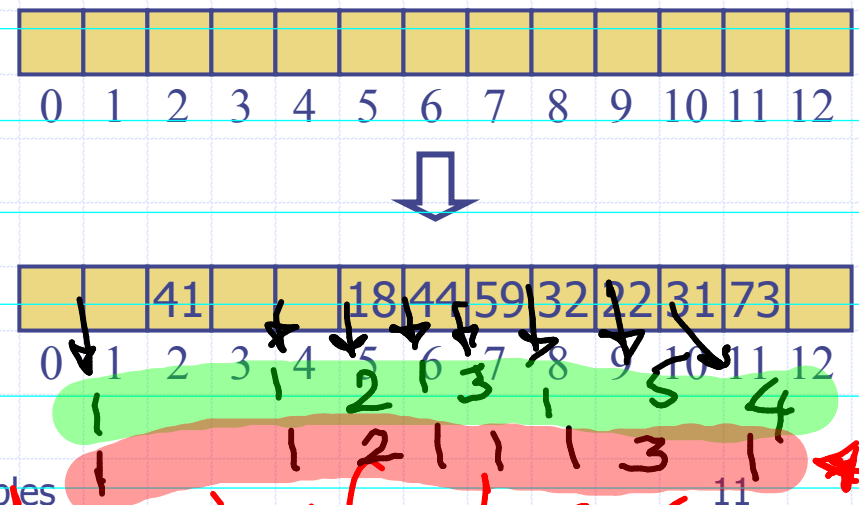
# Linear Probing

(Somewhat overcomes memory wasted in pointers stored in linked lists)

- ◆ **Open addressing**: the colliding item is placed in a different cell of the table
- ◆ **Linear probing** handles collisions by placing the colliding item in the next (circularly) available table cell
- ◆ Each table cell inspected is referred to as a "probe"
- ◆ Colliding items lump together, causing future collisions to cause a longer sequence of probes

## Example:

- $h(x) = x \bmod 13$
- Insert keys 18, 41, 22, 44, 59, 32, 31, 73, in this order



Exercise: Compare total cost of linear probing with linked list for this problem