

# Extended Introduction to Computer Science

## CS1001.py

### Lecture 18a:

## Hash Functions and Hash Tables (cont.)

Instructors: Jonathan Berant, Amir Rubinstein

Teaching Assistants: Michal Kleinbort, Noam Parzanchevski, Ben Bogin

School of Computer Science

Tel-Aviv University

Spring Semester 2019

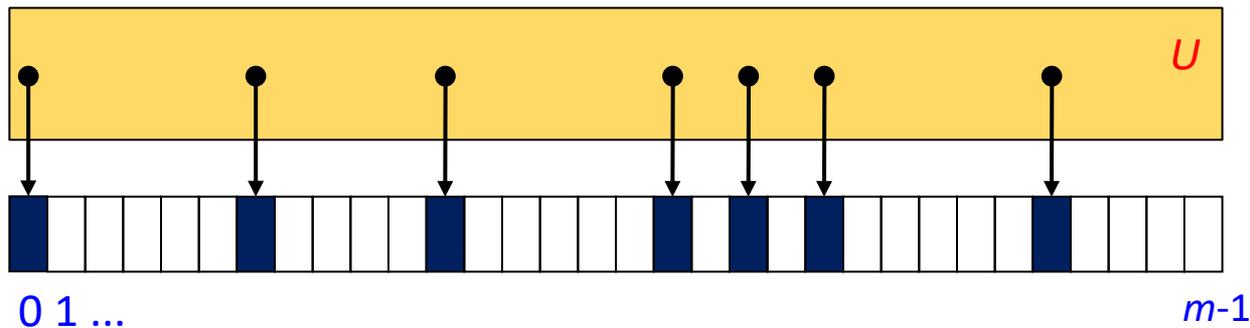
<http://tau-cs1001-py.wikidot.com>

# Goal of Hash Tables

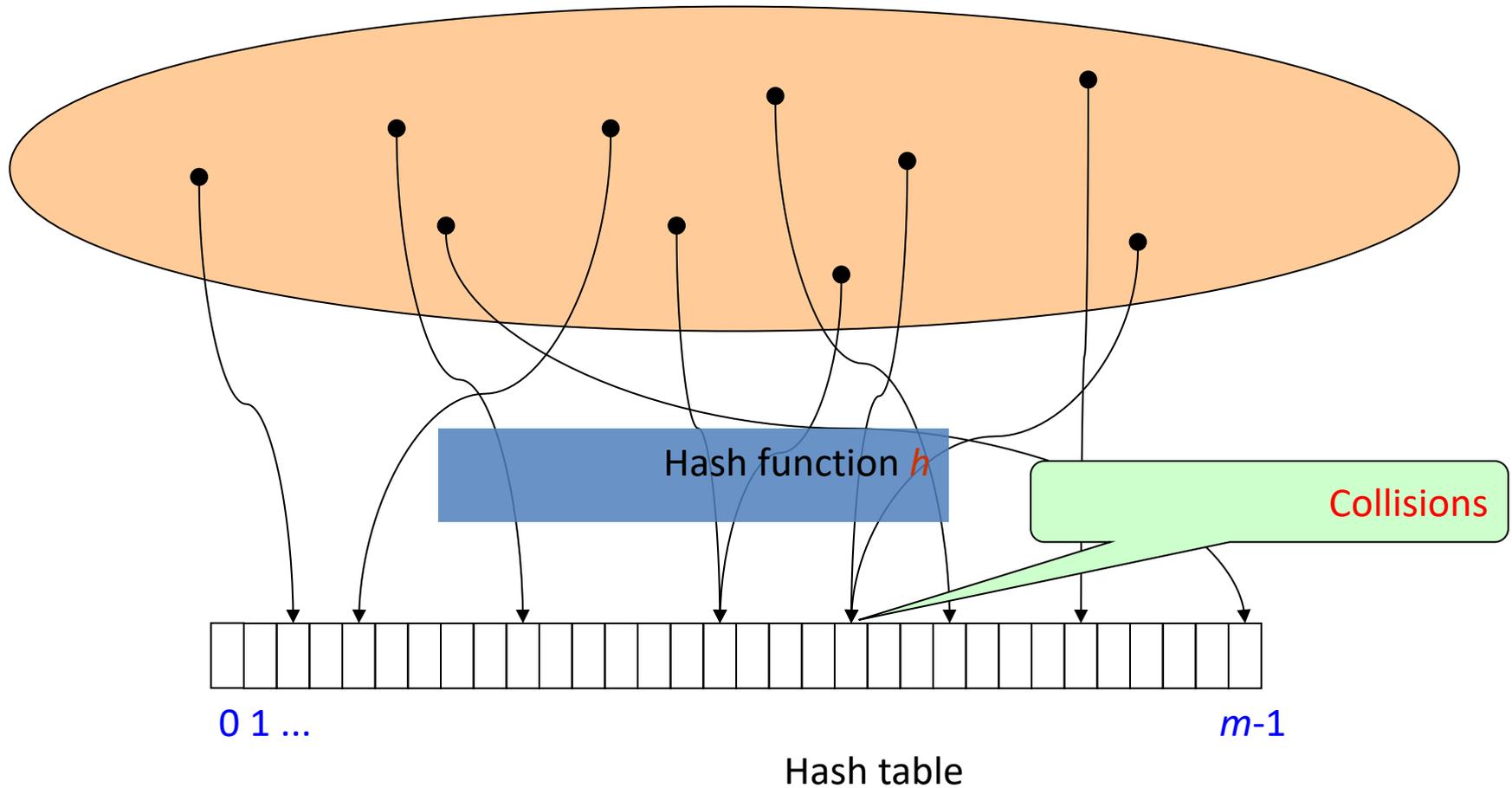
- A data structure  $D$  that allows to efficiently:
  - Insert a key  $k$  into  $D$
  - Find a key  $k$  in  $D$
  - Delete a key  $k$  from  $D$
- We can use balanced binary search trees to obtain  $O(\log n)$  time complexity
- Our goal is to get better average time complexity

# Goal of Hash Tables

- Suppose all possible keys in the universe  $U$  can be stored in an array of size  $m=|U|$
- We can use an array where keys are indices (direct addressing)



# Large Universe ( $|U| \gg n$ )

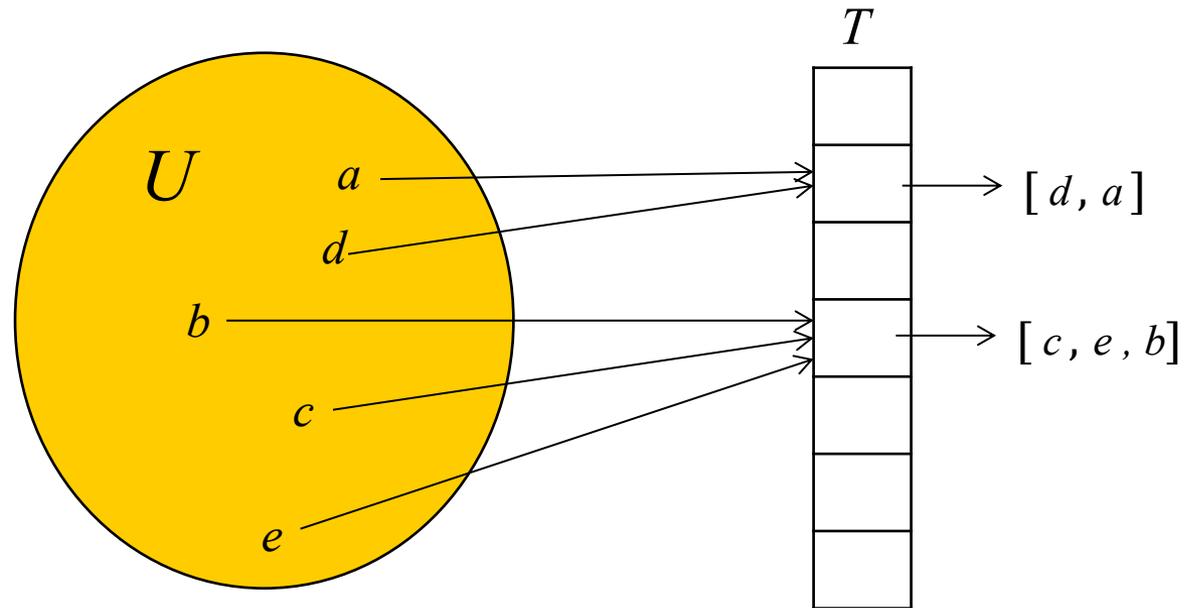


# Hash Function

1. Should be easy to compute
2. Behave like a “random” function (but it is deterministic)

We assume an ideal hash function that distributes keys **uniformly** and **independently**

# Dealing with collisions: **Chaining** (Reminder)



- Each **cell** in the table will contain a **list**
- Search: search in the list  $h(k)$
- Insert: append to the list  $h(k)$
- Delete: remove from the list  $h(k)$

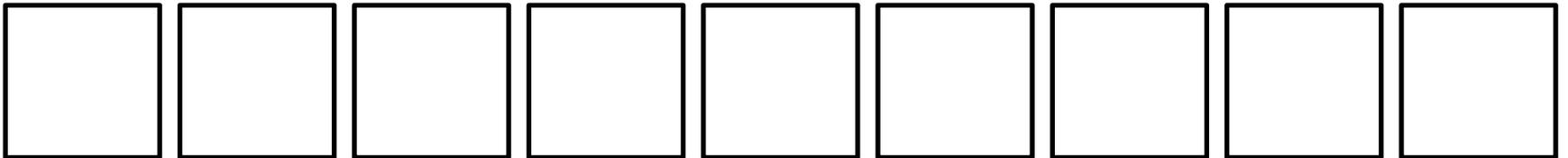
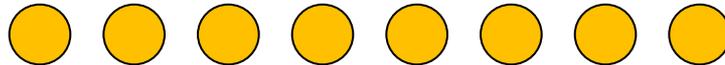
# Hashing with chaining with a random hash function

## Balls in Bins

Throw  $n$  balls randomly into  $m$  bins

# Balls in Bins

Throw  $n$  balls randomly into  $m$  bins



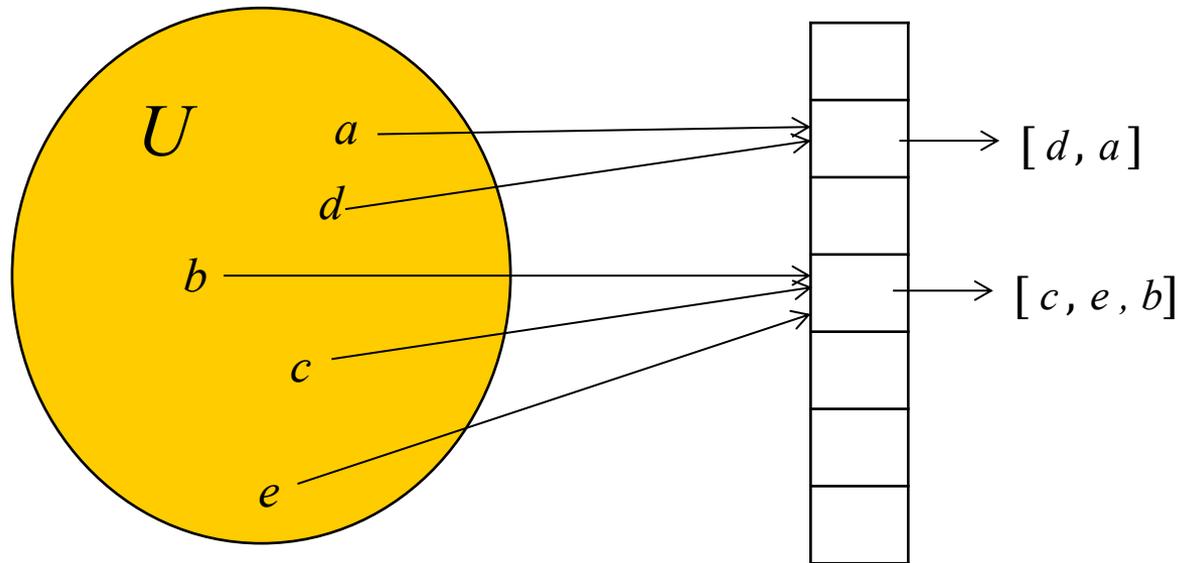
All throws are uniform and independent

# Balls in Bins

Throw  $n$  balls randomly into  $m$  bins

Average number of balls  
in each bin is  $n/m$

# Chaining – Time Complexity $T$



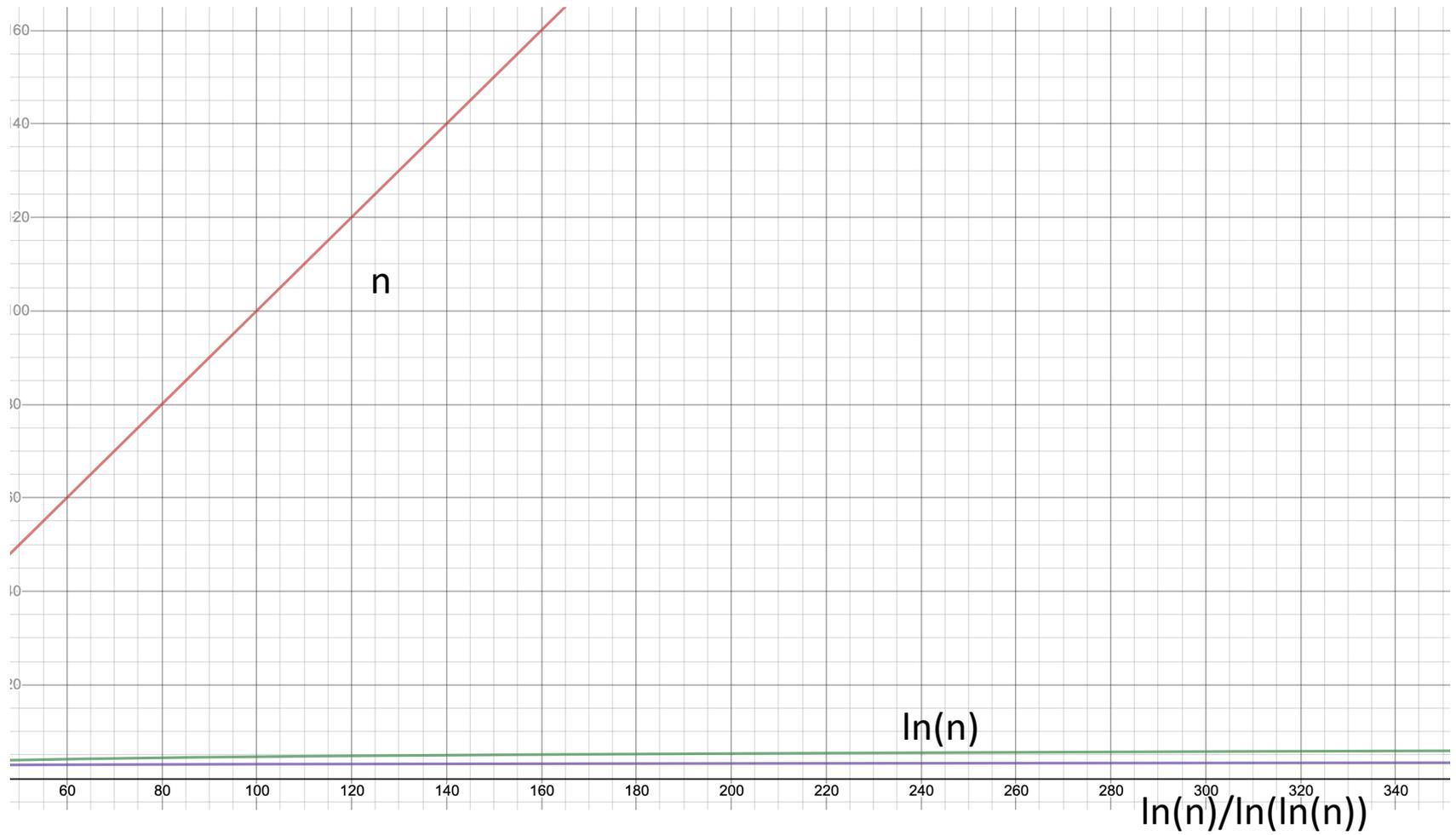
- The **average** length of a chain is  $n/m$ . This is called the "**load factor**", denoted  $\alpha$ .
- if we choose  $m$  such that  $n = O(m)$ , then  $\alpha = O(1)$ .
- This requires some **estimation** of the number of elements we **expect** to be in the table, or a mechanism to dynamically update the table size
- Average time complexity when  $\alpha = O(1)$  is constant. Worst case is always  $O(n)$

# Collision Size – for reference only

- Let  $|\mathcal{K}| = n$  and  $|T| = m$ .
- The expected maximal capacity in a single slot is known for the following cases (proof omitted):

load factor	condition	expected maximal capacity is a single slot
sublinear	$n < \sqrt{m}$	1 (=no collisions)
	$n = m^{1-\epsilon}$ , $0 < \epsilon < 1/2$	$O(1/\epsilon)$
linear	$n = m$	$\frac{\ln(n)}{\ln \ln(n)}$
superlinear	$n > m$	$\frac{n}{m} + \frac{\ln(n)}{\ln \ln(n)}$

- This provides some intuition as to the **rareness of the worst case**.



# Implementation in Python

- Let us implement our own class `Hashtable` in Python now.
- We will use `chaining` to resolve collisions.
- We will demonstrate its usage with simple strings first. Later on we will show another example with class `Student`.

# Initializing the Hash Table

```
class Hashtable:

    def __init__(self, m, hash_func=hash):
        """ initial hash table, m empty entries """
        self.table = [ [] for i in range(m) ]
        self.hash_mod = lambda key: hash_func(key) % m

    def __repr__(self):
        return "".join([str(i) + " " + str(self.table[i]) + "\n"
                        for i in range(len(self.table))])
```

# Initializing the Hash Table

```
>>> ht = Hashtable (11)
```

```
>>> ht
```

```
0 []
```

```
1 []
```

```
2 []
```

```
3 []
```

```
4 []
```

```
5 []
```

```
6 []
```

```
7 []
```

```
8 []
```

```
9 []
```

```
10 []
```

Since our table is a list of lists, and lists are **mutable**, we should be **careful** even when initializing the list.

# Initializing the Hash Table: a **Bogus** Code

Consider the following alternative initialization:

```
class Hashtable:
    def __init__(self, m, hash_func=hash):
        """ initial hash table, m empty entries """
        self.table = [[]]*m
```

```
>>> ht = Hashtable(11)
>>> ht.table[0].append(5)
>>> ht
0 [5]
1 [5]
...
>>> ht.table[0] == ht.table[1]
True
>>> ht.table[0] is ht.table[1]
True
```

The entries produced by this **bogus** `__init__` are **identical**.  
Therefore, mutating one mutate all of them:

# Initializing the Hash Table, cont.

But this one will work fine, as different copies of the empty list will be created:

```
class Hashtable:

    def __init__(self, m, hash_func=hash):
        """ initial hash table, m empty entries """
        empty = []
        self.table = [ list(empty) for i in range(m) ]
```

# Dictionary Operations: Python Code

```
class Hashtable:
...

def find(self, item):
    """ returns True if item in hashtable, False otherwise """
    i = self.hash_mod(item)
    chain = self.table[i]
    if item in chain: # a hidden loop
        return True
    else:
        return False

def insert(self, item):
    """ insert an item into table, if not there """
    i = self.hash_mod(item)
    chain = self.table[i]
    if item not in chain: # a hidden loop
        chain.append(item)
```

return item in chain

# Example: A **Very** Small Table

( $n = 14$ ,  $m = 7$ )

In the following slides, there are executions that construct a hash table with  $m = 7$  entries. We'll insert  $n = 14$  string records in it and check how insertions are distributed, and in particular what the maximum number of collisions per cell is.

Our hash table will be a **list** with  $m = 7$  entries. Each entry will contain a list with a **variable length**. Initially, each entry of the hash table is an **empty list**.

# Example: A **Very** Small Table

## (n = 14, m = 7)

```
>>> tribes = ['Reuben', 'Simeon', 'Levi', 'Judah', 'Dan',  
             'Naphtali', 'Gad', 'Asher', 'Issachar', 'Zebulun', 'Benjamin',  
             'Joseph', 'Ephraim', 'Manasse']
```

```
>>> ht = Hashtable(7)
```

```
>>> for name in tribes:  
    ht.insert(name)
```

```
>>> ht #calls __repr__  
      (next slide)
```

# Example: A **Very** Small Table

(n = 14, m = 7)

```
>>> ht #calls __repr__
0 []
1 ['Reuben', 'Judah', 'Dan']
2 ['Naphtali']
3 ['Gad', 'Ephraim']
4 ['Levi']
5 ['Issachar', 'Zebulun']
6 ['Simeon', 'Asher', 'Benjamin', 'Joseph', 'Manasse']
```

# Example: A slightly larger table (n = 14, m = 21)

```
>>> tribes = ['Reuben', 'Simeon', 'Levi', 'Judah', 'Dan',  
             'Naphtali', 'Gad', 'Asher', 'Issachar', 'Zebulun', 'Benjamin',  
             'Joseph', 'Ephraim', 'Manasse']
```

```
>>> ht = Hashtable(21)
```

```
>>> for name in tribes:  
    ht.insert(name)
```

```
>>> ht #calls __repr__  
      (next slide)
```

# Example: A slightly larger table (n = 14, m = 21)

```
>>> ht #calls __repr__
0 []
1 []
2 []
3 ['Ephraim']
4 []
5 ['Issachar']
6 ['Benjamin']
7 []
8 ['Judah']
9 ['Naphtali']
10 []
11 []
12 ['Zebulun']
13 ['Manasse']
14 []
15 ['Reuben', 'Dan']
16 []
17 ['Gad']
18 ['Levi']
19 []
20 ['Simeon', 'Asher', 'Joseph']
```

# Hashing random seed

- If you run this code yourself, you will probably encounter results that are **different** from those in the **last few slides**.
- This is because when **IDLE starts**, it **randomly** generates a number called **seed**, which is used to compute the built-in **hash** function
- This is intended to provide **protection** against **denial-of-service attacks** caused by carefully-chosen inputs designed to **collide**. Such attacks exploit the **worst case** performance of using **hash**.
- Of course, as long as you work under the **same instance of IDLE**, **hash** is **consistent**. But re-opening IDLE will probably break this consistency.

# Hashing and User-defined Classes

- So far we used our `Hashtable` class to store Python's built-in types such as `str` and `int`.
- We will now use class `Hashtable` on our own `class Student`.
- As we will see, this will raise certain issues, which we will solve.

# The Student Class (reminder)

```
class Student:
```

```
    def __init__(self, name, surname, ID):
```

```
        self.name = name
```

```
        self.surname = surname
```

```
        self.id = ID
```

```
        self.grades = dict()
```

```
    def __repr__(self): #must return a string
```

```
        return "<" + self.name + ", " + str(self.id) + ">"
```

```
    def update_grade(self, course, grade):
```

```
        self.grades[course] = grade
```

```
    def avg(self):
```

```
        s = sum([self.grades[course] for course in self.grades])
```

```
        return s / len(self.grades)
```

# Hashing Students

```
>>> st1 = Student("Grace", "Hopper", 123456789)
>>> st2 = Student("Grace", "Hopper", 123456789)
```

```
>>> st1
<Grace, 123456789>
>>> st2
<Grace, 123456789>
```

```
>>> hash(st1)
-9223372036851698786
>>> hash(st2)
3077117
```



## From Wikipedia:

**Grace Brewster Murray Hopper** (1906 –1992), was an American computer scientist and United States Navy Rear Admiral. She was one of the first programmers of the [Harvard Mark I](#) computer in 1944, invented the first compiler for a computer programming language, and was one of those who popularized the idea of machine-independent programming languages which led to the development of [COBOL](#), one of the first [high-level programming languages](#).

- This should not be a surprise to you: by **default**, Python uses the **memory address** of an object to compute the value of **hash** on it.

# The `__hash__` method

- We will add to `class Student` the special method `__hash__`.
- It defines the result of calling Python's `hash` on an object of this class.

```
class Student:
```

```
...
```

```
    def __hash__(self): #so we can use hash(st) on a student  
        return hash(self.id) #assume id is a unique identifier
```

- Notes:
  - 1) `__hash__` of `Student` class calls `__hash__` of `int` class
  - 2) We used **merely the id** to compute a student's hash, under the assumption that it is **unique**. We could have used **more fields of a Student object**.

# Hashing Students – almost done

```
>>> st1 = Student("Grace", "Hopper", 123456789)
```

```
>>> st2 = Student("Grace", "Hopper", 123456789)
```

```
>>> hash(st1) == hash(st2) == hash(st1.id)
```

```
True 😊
```

# Hashing Students – almost done

- Can you explain why the following search **fails**?

```
>>> st1 = Student("Grace", "Hopper", 123456789)
>>> st2 = Student("Grace", "Hopper", 123456789)
```

```
>>> ht = Hashtable(7)
```

```
>>> ht.insert(st1)
```

```
>>> ht
```

```
0 []
```

```
1 [<Grace, 123456789>]
```

```
2 []
```

```
3 []
```

```
4 []
```

```
5 []
```

```
6 []
```

```
>>> ht.find(st2)
```

```
False 😞
```

# hashing requires `__eq__` as well

- Indeed, **no much point** in having `__hash__` without `__eq__`, for comparing elements (within a chain inside a table's index).

```
class Student :
```

```
...
```

```
def __eq__(self, other):  
    return self.name == other.name and \  
           self.surname == other.surname and \  
           self.id == other.id
```

```
>>> ht.find(st2) # recall st2 holds same data as st1  
True 😊
```

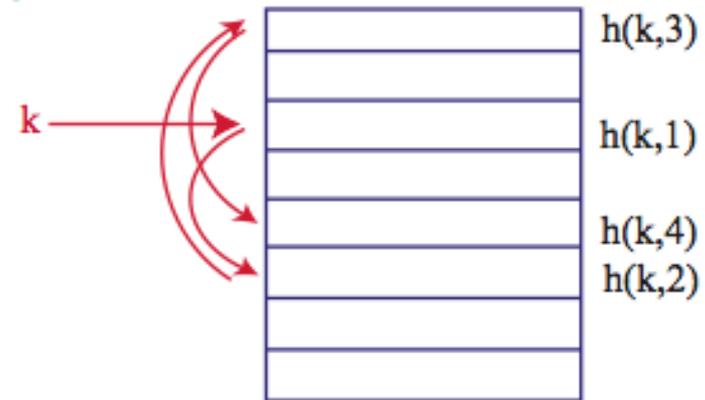
# Two Approaches for Dealing with Collisions

- 1) **Chaining** – explained and implemented this ✓
- 2) **Open addressing** – we will briefly discuss it now

# Two Approaches for Dealing with Collisions:

## (2) Open Addressing

- In open addressing, each slot in the hash table contains **at most one** item. This obviously implies that  **$n$  cannot be larger than  $m$** .
- Furthermore, an item will typically not stay statically in the slot where it "tried" to enter, or where it was placed initially. Instead, it may be moved a few times around.



- Open addressing is important in **hardware applications** where devices have many slots but each can only store one item (e.g. fast **switches** and high capacity **routers** ). It is also used in **python dictionaries and sets**.
- There are many approaches to open addressing. We will describe a fairly recent one, termed **cuckoo hashing** (Pagh and Rodler, 2001).

# Cuckoo Hashing: Motivation

- We saw that if  $n \leq m$ , hashing with chaining guarantees that insertion, deletion, and find are carried out in expected time  $O(1)$  per operation, and **with high probability** (probability is over choices of inputs)  $O(\log n / \log \log n)$  per operation. (The worst case time is  $O(n)$  per operation.)
- In certain scenarios (e.g. fast routers in large internet nodes) we want **find** to run in  $O(1)$  **worst-case** time.
- **Cuckoo hashing** is one way to achieve this, but there are two prices to pay:
  - 1) Instead of  $n \leq m$ , we require  $n \leq m/2$ .
    - That is, we pay a price in terms of memory
  - 2) **insert** may take **somewhat longer time**.

# Cuckoo Hashing

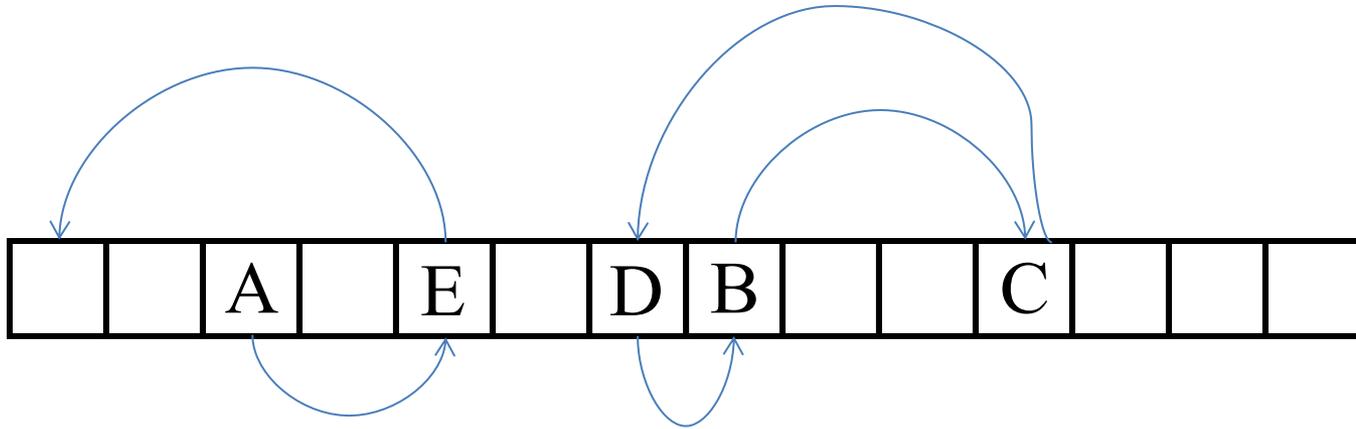
- **Cuckoo hashing** uses two distinct hash functions,  $h_1$  and  $h_2$  (improved versions use four, six, or eight, but the idea is the same).
- Each key,  $k$ , has **two potential slots** in the hash table,  $h_1(k)$  and  $h_2(k)$ . If we search for  $k$ , all we have to do is look for it in these two locations (no chains here -- at most **one item** per slot).
- It is slightly more involved to **insert** a record whose key is  $k$ .

# Cuckoo Hashing

It is slightly more involved to **insert** a record whose key is  $k$ .

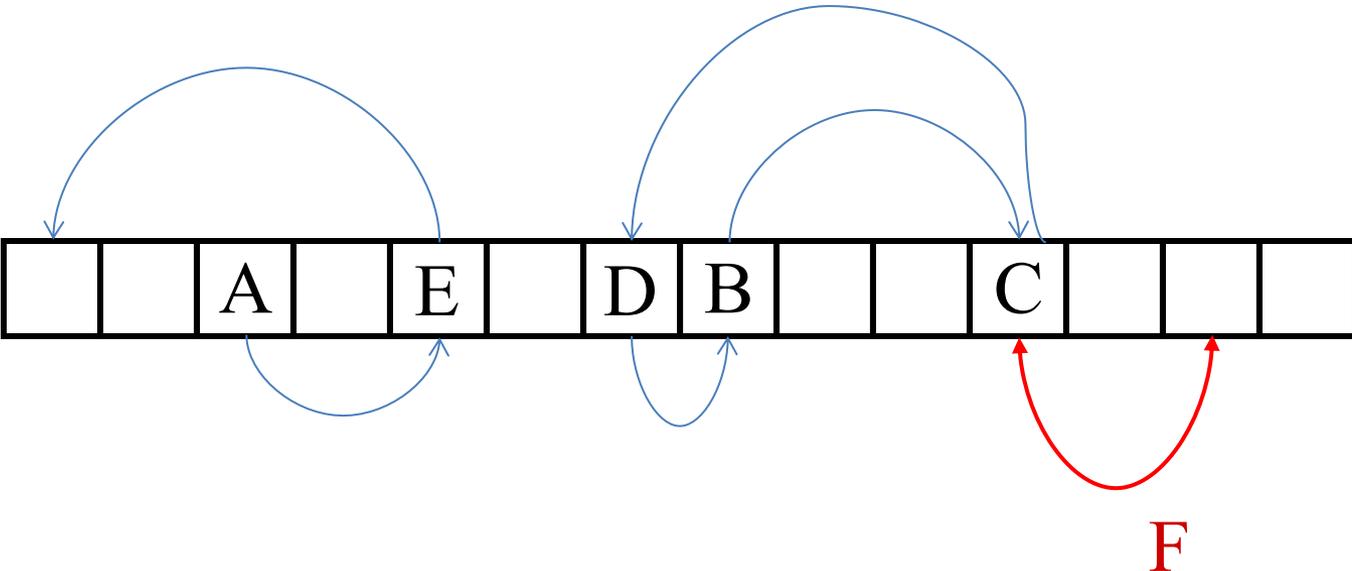
- If any of the two slots,  $h_1(k)$  or  $h_2(k)$  is empty,  $k$  is inserted there.
- If both slots are full, pick one of the two occupants, say  $x$ . Place  $k$  in  $x$ 's current slot.
- Assume this was location  $h_1(x)$ . Place  $x$  in its other slot,  $h_2(x)$ .
- If that slot was empty, we are done.
- Otherwise, the slot is occupied by some  $y$ . Place this  $y$  in its other slot, potentially kicking its present occupant, etc.,etc., until we find an empty slot.

# Cuckoo Hashing: Examples



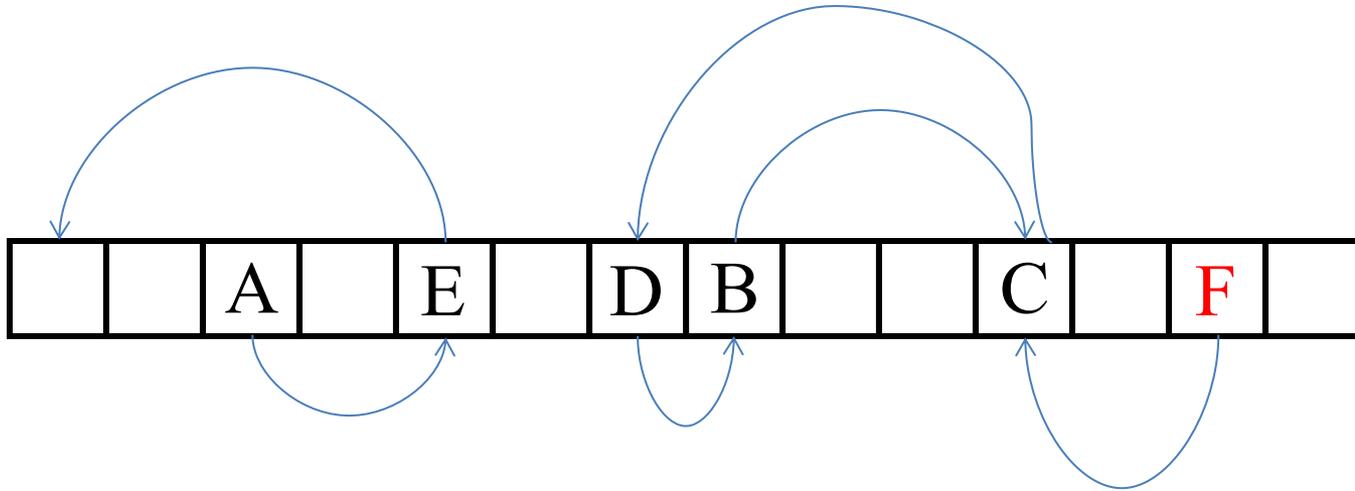
= The other potential slot for an item

# Cuckoo Hashing: Examples



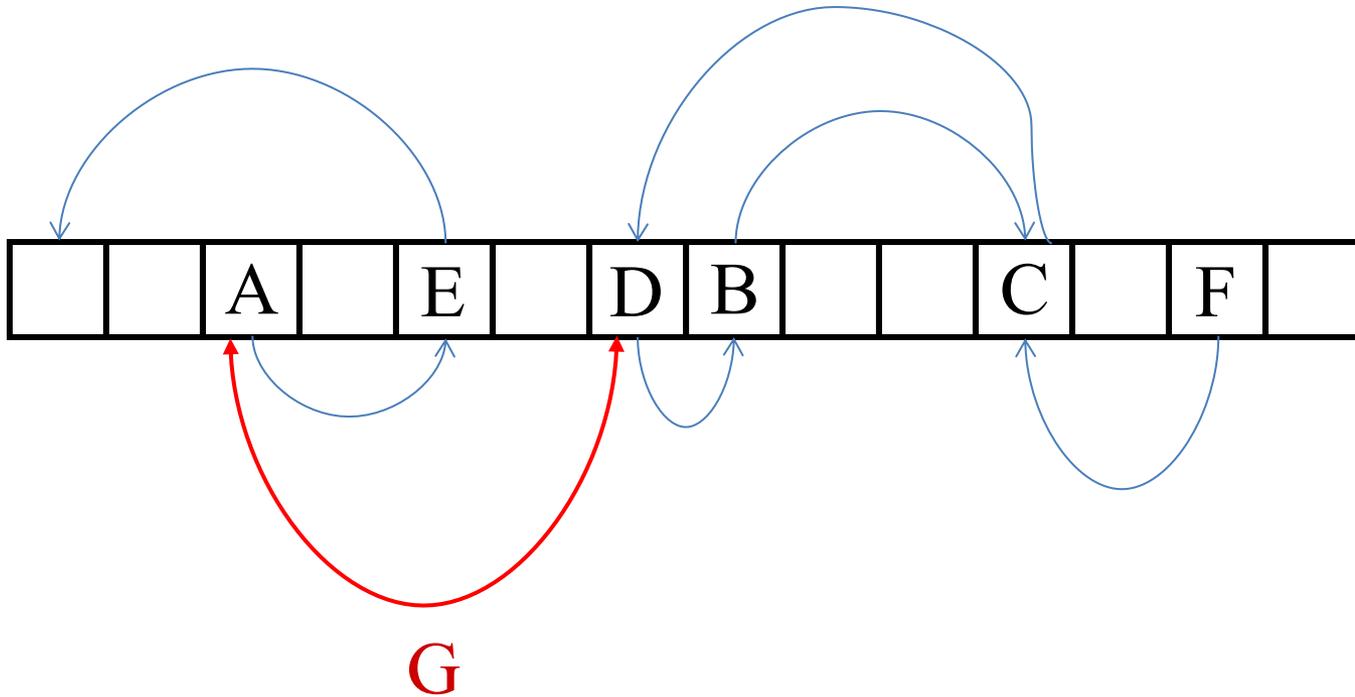
= The other potential slot for an item

# Cuckoo Hashing: Examples



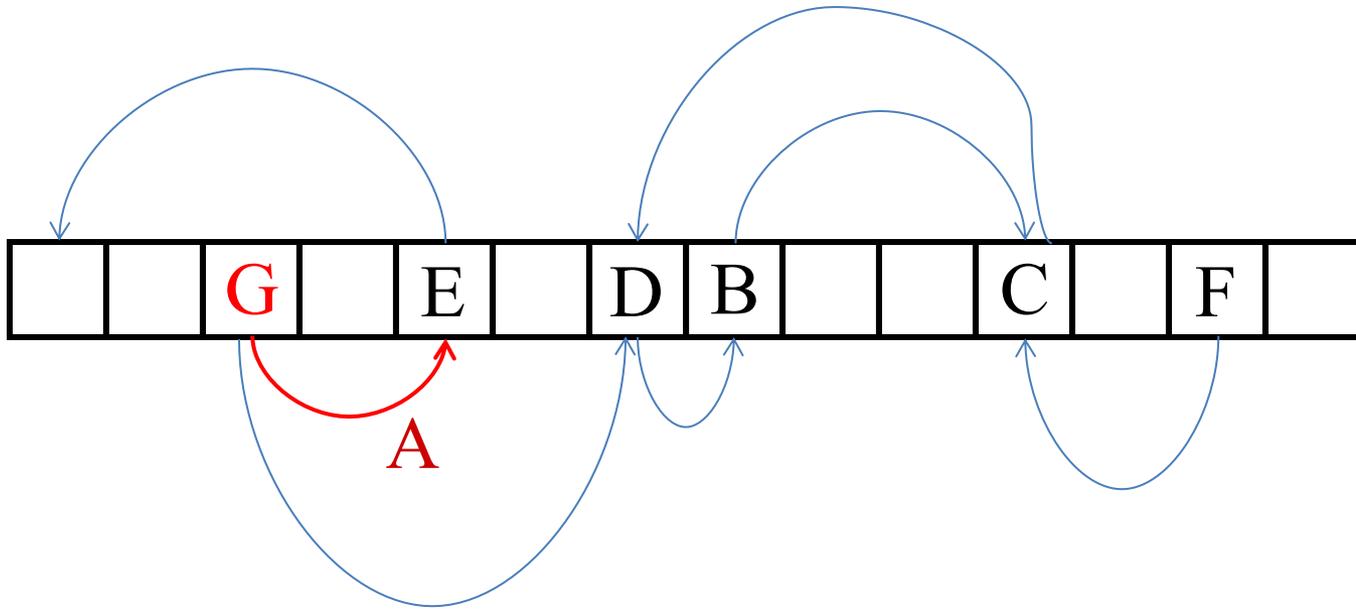
= The other potential slot for an item

# Cuckoo Hashing: Examples



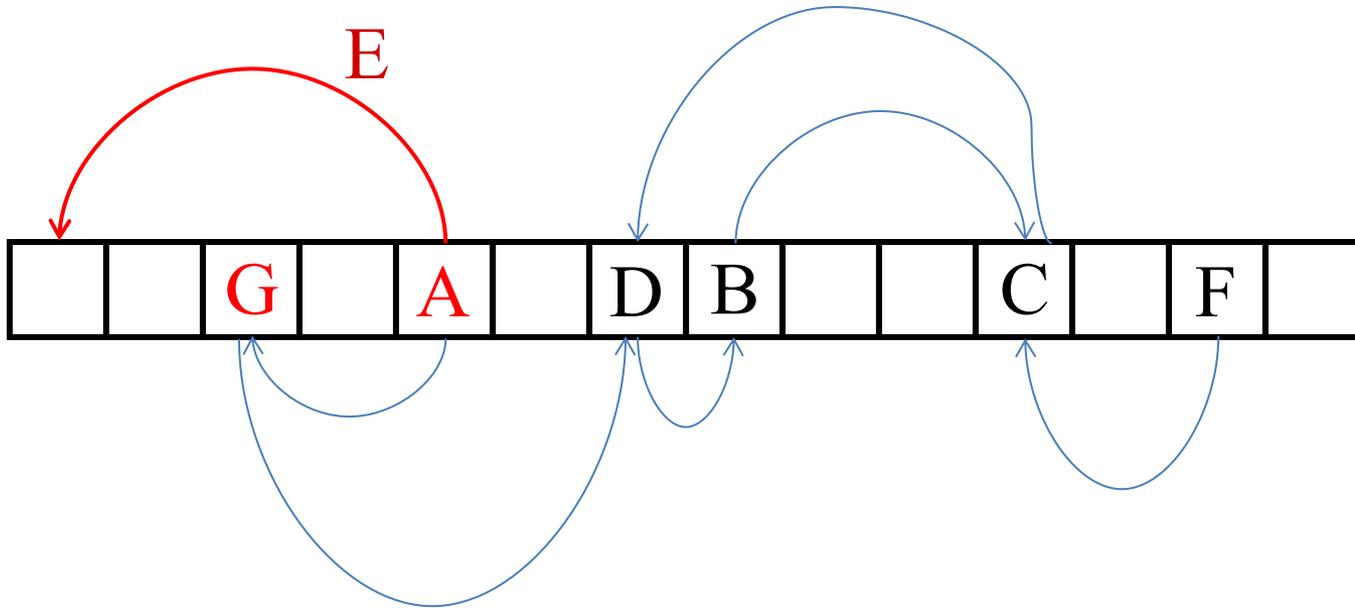
= The other potential slot for an item

# Cuckoo Hashing: Examples



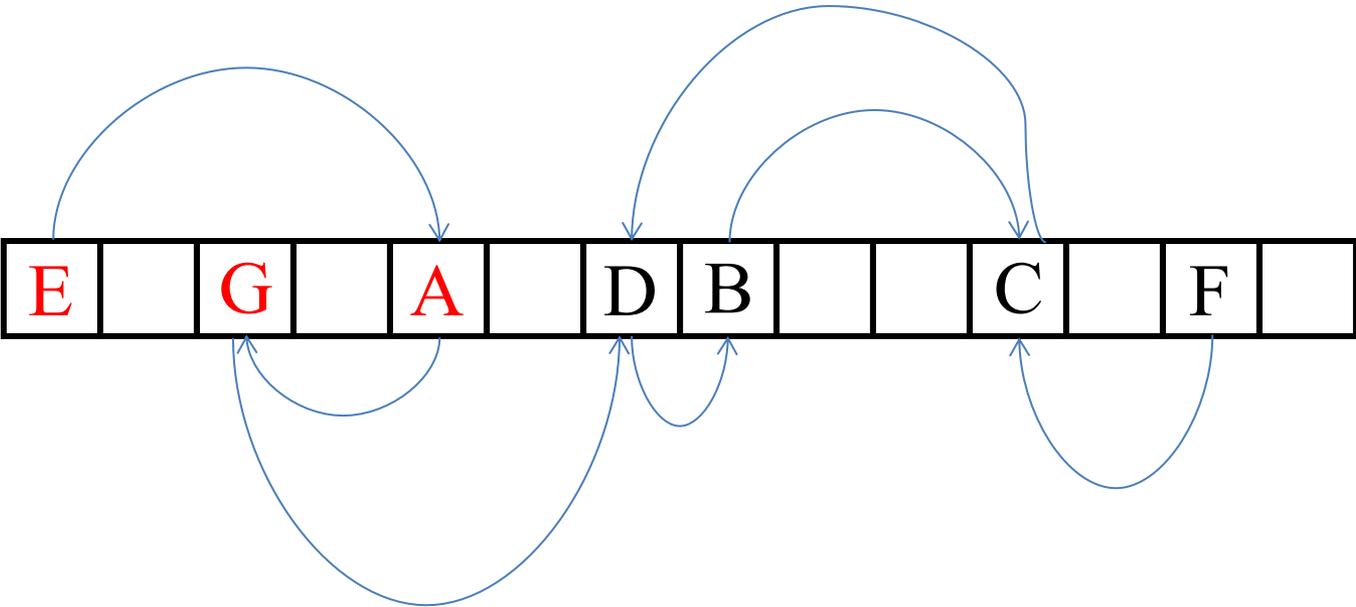
= The other potential slot for an item

# Cuckoo Hashing: Examples



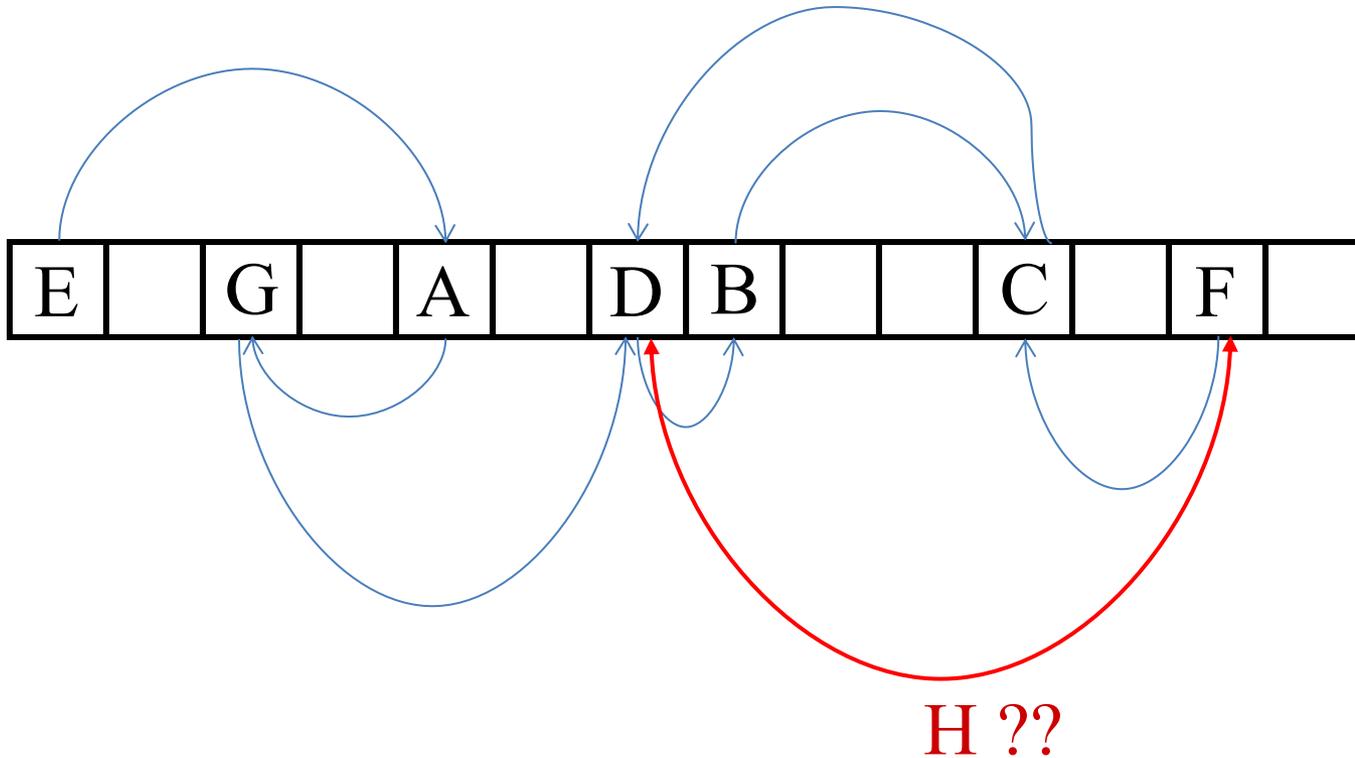
= The other potential slot for an item

# Cuckoo Hashing: Examples



= The other potential slot for an item

# Cuckoo Hashing: Examples



= The other potential slot for an item

# Cuckoo Hashing - Deadlocks

- In the last example, we have reached a **cycle**, and we are in a non ending loop. This is called a **deadlock**.
- The union of the potential locations of **5 items** (B, C, D, F, H) is just **4 slots**.
- This obviously is very bad news for our cuckoo hashing.
- Notice that this is not a very likely event. With very high probability, the 10 potential locations ( $10=5 \cdot 2$ ) will attain **more** than just 4 distinct values (which is why we got stuck in the last example).

# Cuckoo Hashing – Solving Deadlocks

- Another possible problem is that there will be no cycle, but the path leading to the successful insertions will be very **long**.
- Fortunately, such unfortunate cases occur with very low probability when the load factor , i.e.  $n/m$ , is **sufficiently low**. The common recommendation for two hash functions ,  $h_1(\cdot)$ ,  $h_2(\cdot)$  , is to have  $n/m < 1/2$  . (More hash functions enable a higher load factor).
- A theoretical solution: In case of failure (or very long path), **rehash** using “**fresh hash functions**.”
- A more practical solution: Maintain a very small **excess zone** (e.g. 32 excess slots for a hash table with  $m=10000$  slots) and place items “causing trouble” there. If regular search (applying  $h_1(x)$ ,  $h_2(x)$  ) fails, search the excess zone as well.

# Cuckoo Hashing in the Real World

- The load factor has to be smaller than 1. Yet a small load factor, say  $n/m < 1/2$ , is a **waste of memory**.
- In high performance routers, for example, most operations (including the hashing) are done **in silico, by the hardware**. The critical resource is memory area within the chip. Low load factor means wasted area.
- Instead of just 2 hash functions, **4 to 8** hash functions are utilized. This allows to increase the load factor to  $n/m = 3/4$  or even  $n/m = 7/8$ .
- Suppose we use **4** hash functions,  $h_1()$ ,  $h_2()$ ,  $h_3()$ ,  $h_4()$ . Given an element,  $x$ , that we wish to insert, we first check if any of the four locations  $h_1(x)$ ,  $h_2(x)$ ,  $h_3(x)$ ,  $h_4(x)$  is free.

# Cuckoo Hashing in the Real World, cont.

- If these 4 locations are all taken, let  $a, b, c, d$  be the four elements in the above mentioned locations, respectively.
- Look, for example, at  $a$ . If one of the other 3 locations among  $h_1(a), h_2(a), h_3(a), h_4(a)$  is free, we move  $a$  there, and put  $x$  in its place. If not, we do the same with respect to  $b$ , then  $c$ , then  $d$ .
- If all these are taken ( $4+4\cdot 3=16$  different locations, typically), we go one more level down this search tree ( $12\cdot 3 = 36$  additional locations, typically).
- If all these are taken, we give up on  $x$  and put it in the garbage bin (“excess zone” table).
- With very high probability, the small excess zone does not fill up. After removing elements from the table, we could try re-inserting such  $x$  to the hash table.

# Designing Distinct Hash Functions

- Recall that the goal of designing a hash functions is that they map most sets of keys such that the maximal number of collisions is **small**.
- When having more than one hash function, we have the additional goal that the different functions map same keys approximately **independently**. In Python, we could try variants of good ole hash.

For example:

```
def hash0(x):  
    return hash("0" + str(x))  
def hash1(x):  
    return hash("1" + str(x))  
def hash2(x):  
    return hash(str(x) + "2")  
def hash3(x):  
    return hash(str(x) + "3")
```

# Designing Distinct Hash Functions

A reminder concerning str (mapping objects to representing strings):

```
>>> [str(i) for i in range(10,20)]
['10', '11', '12', '13', '14', '15', '16', '17', '18', '19']
>>> str(2.2)
'2.2'
>>> str("2.2")
'2.2'
```

And now applying the four functions on a small domain:

```
>>> for f in (hash0,hash1,hash2,hash3):
    print([f(i) %23 for i in range(10,20)])
[5, 1, 12, 21, 5, 10, 7, 0, 11, 17]
[17, 18, 18, 8, 19, 4, 2, 1, 22, 11]
[16, 18, 19, 3, 8, 6, 13, 15, 3, 20]
[17, 8, 18, 0, 19, 7, 1, 14, 8, 14]
```

Random? Independent? Mixing well? You be the judges.

# Data Structures Wrap Up:

## Search

- Search has always been a **central computational task**. The emergence and the popularization of the world wide web has literally created a **universe of data**, and with it the need to pinpoint information in this universe.
- Various **search engines** have emerged, to cope with this **big data** challenge. They constantly collect data on the web, organize it, and store it in sophisticated data structures that support efficient (very fast) access, resilience to failures, frequent updates, including deletions, etc. etc.

# Data Structures Wrap Up:

## Dictionaries in Computer Science

- In computer science, a **dictionary** is an abstract data type (ADT) supporting efficient **insert**, **delete**, and **search** operations.
- This is an abstract notion, and should **not** be confused with Python's **class dict**, although **class dict** can be thought of as an implementation of the abstract data type **dictionary** (with elements that are pairs key:value).
- There are **two variations** of this ADT, according to the type of the elements stored in the data structure.
  - pairs **key: value**, or
  - Just **keys**

In any case, we assume all keys are **unique** (different items have different keys).

- We saw several possible implementation for dictionaries:
  - Lists (as arrays or as linked lists, sorted or unsorted)
  - Balanced search trees
  - Hash tables

# Data Structures Wrap Up:

## Hash Functions and Tables

- Hash functions map large domains to smaller ranges.
- Example:  
$$h : \{0,1,\dots,p^2\} \rightarrow \{0,1,\dots,p-1\},$$

defined by  $h(x) = a \cdot x + b \pmod{p}$ .
- Hash tables are extensively used for searching.
- If the range is smaller than the domain, we cannot avoid **collisions** ( $x \neq y$  with  $h(x) = h(y)$ ). For example, in the example above, if  $x_1 = x_2 \pmod{p}$  then  $h(x_1) = h(x_2)$ .
- If the domain size is larger than **the square root** of the range size, there will be **collisions** with high probability.
- A good hash function should create few collisions for most subsets of the domain ("few" is relative to size of subset).

# Data Structures Wrap Up: Using Hash Functions and Tables

- We explained chaining as a way to resolve collisions.
- We also studied the paradigm of cuckoo hashing , using two hash functions  $h_1()$ ,  $h_2()$  (or four, or eight).  
Cuckoo hashing is aimed at a constant time find operation (**worst-case**), at the cost of a slightly longer insert operation.
- In the data structures course, you will see additional collisions resolution means, such as double hashing, etc.
- Python **sets** and **dictionaries** use hash tables, thus searching an element in a set / dict takes  $O(1)$  time on average. Collisions are solved using open addressing, in a more sophisticated manner. In addition, the size of the hash table is dynamic.