

# SE274 Data Structure

## Lecture 6: Maps, Hash Tables, Skip Lists – Part 1

Apr 06, 2020

Instructor: Sunjun Kim

Information&Communication Engineering, DGIST

# Maps and Dictionaries



# Maps



- A **map** is a searchable collection of items that are key-value pairs
- The main operations of a map are for searching, inserting, and deleting items
- Multiple items with the same key are **not** allowed
- Applications:
  - address book
  - student-record database

# Dictionaries

- Python's **dict** class is arguably the most significant data structure in the language.
  - It represents an abstraction known as a **dictionary** in which unique **keys** are mapped to associated **values**.
- Here, we use the term “dictionary” when specifically discussing Python's dict class, and the term “map” when discussing the more general notion of the abstract data type.

# The Map ADT (Using dict Syntax)



- $M[k]$ : Return the value  $v$  associated with key  $k$  in map  $M$ , if one exists; otherwise raise a `KeyError`. In Python, this is implemented with the special method `__getitem__`.
- $M[k] = v$ : Associate value  $v$  with key  $k$  in map  $M$ , replacing the existing value if the map already contains an item with key equal to  $k$ . In Python, this is implemented with the special method `__setitem__`.
- `del M[k]`: Remove from map  $M$  the item with key equal to  $k$ ; if  $M$  has no such item, then raise a `KeyError`. In Python, this is implemented with the special method `__delitem__`.
- `len(M)`: Return the number of items in map  $M$ . In Python, this is implemented with the special method `__len__`.
- `iter(M)`: The default iteration for a map generates a sequence of *keys* in the map. In Python, this is implemented with the special method `__iter__`, and it allows loops of the form, **for  $k$  in  $M$ .**

# More Map Operations

`k in M`: Return `True` if the map contains an item with key `k`. In Python, this is implemented with the special `__contains__` method.

`M.get(k, d=None)`: Return `M[k]` if key `k` exists in the map; otherwise return default value `d`. This provides a form to query `M[k]` without risk of a `KeyError`.

`M.setdefault(k, d)`: If key `k` exists in the map, simply return `M[k]`; if key `k` does not exist, set `M[k] = d` and return that value.

`M.pop(k, d=None)`: Remove the item associated with key `k` from the map and return its associated value `v`. If key `k` is not in the map, return default value `d` (or raise `KeyError` if parameter `d` is `None`).

# A Few More Map Operations

`M.popitem()`: Remove an arbitrary key-value pair from the map, and return a `(k,v)` tuple representing the removed pair. If map is empty, raise a `KeyError`.

`M.clear()`: Remove all key-value pairs from the map.

`M.keys()`: Return a set-like view of all keys of `M`.

`M.values()`: Return a set-like view of all values of `M`.

`M.items()`: Return a set-like view of `(k,v)` tuples for all entries of `M`.

`M.update(M2)`: Assign `M[k] = v` for every `(k,v)` pair in map `M2`.

`M == M2`: Return `True` if maps `M` and `M2` have identical key-value associations.

`M != M2`: Return `True` if maps `M` and `M2` do not have identical key-value associations.

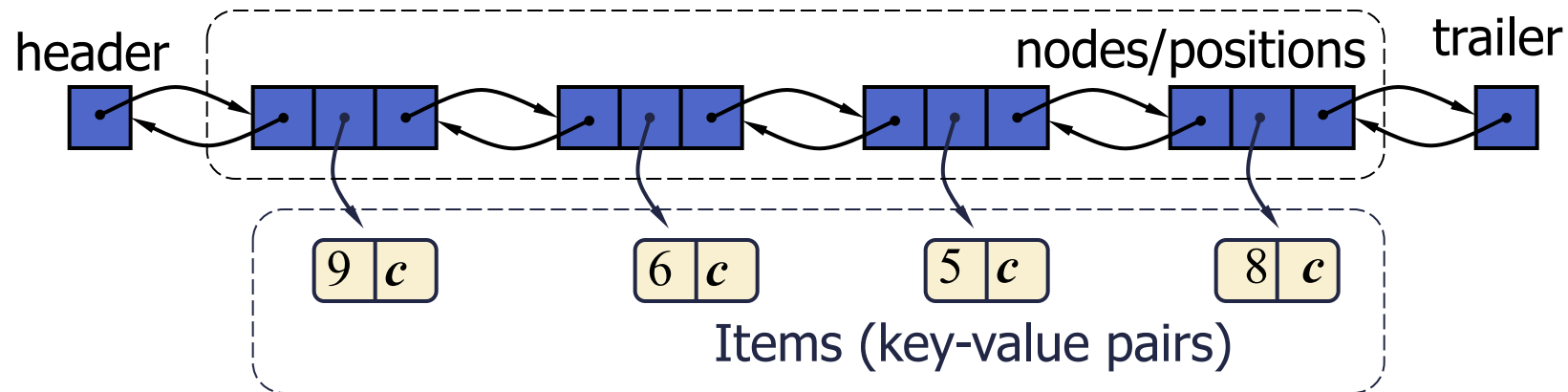
# Example

Operation	Return Value	Map
len(M)	0	{ }
M['K'] = 2	—	{ 'K': 2 }
M['B'] = 4	—	{ 'K': 2, 'B': 4 }
M['U'] = 2	—	{ 'K': 2, 'B': 4, 'U': 2 }
M['V'] = 8	—	{ 'K': 2, 'B': 4, 'U': 2, 'V': 8 }
M['K'] = 9	—	{ 'K': 9, 'B': 4, 'U': 2, 'V': 8 }
M['B']	4	{ 'K': 9, 'B': 4, 'U': 2, 'V': 8 }
M['X']	KeyError	{ 'K': 9, 'B': 4, 'U': 2, 'V': 8 }
M.get('F')	None	{ 'K': 9, 'B': 4, 'U': 2, 'V': 8 }
M.get('F', 5)	5	{ 'K': 9, 'B': 4, 'U': 2, 'V': 8 }
M.get('K', 5)	9	{ 'K': 9, 'B': 4, 'U': 2, 'V': 8 }
len(M)	4	{ 'K': 9, 'B': 4, 'U': 2, 'V': 8 }
del M['V']	—	{ 'K': 9, 'B': 4, 'U': 2 }
M.pop('K')	9	{ 'B': 4, 'U': 2 }
M.keys()	'B', 'U'	{ 'B': 4, 'U': 2 }
M.values()	4, 2	{ 'B': 4, 'U': 2 }
M.items()	('B', 4), ('U', 2)	{ 'B': 4, 'U': 2 }
M.setdefault('B', 1)	4	{ 'B': 4, 'U': 2 }
M.setdefault('A', 1)	1	{ 'A': 1, 'B': 4, 'U': 2 }
M.popitem()	('B', 4)	{ 'A': 1, 'U': 2 }

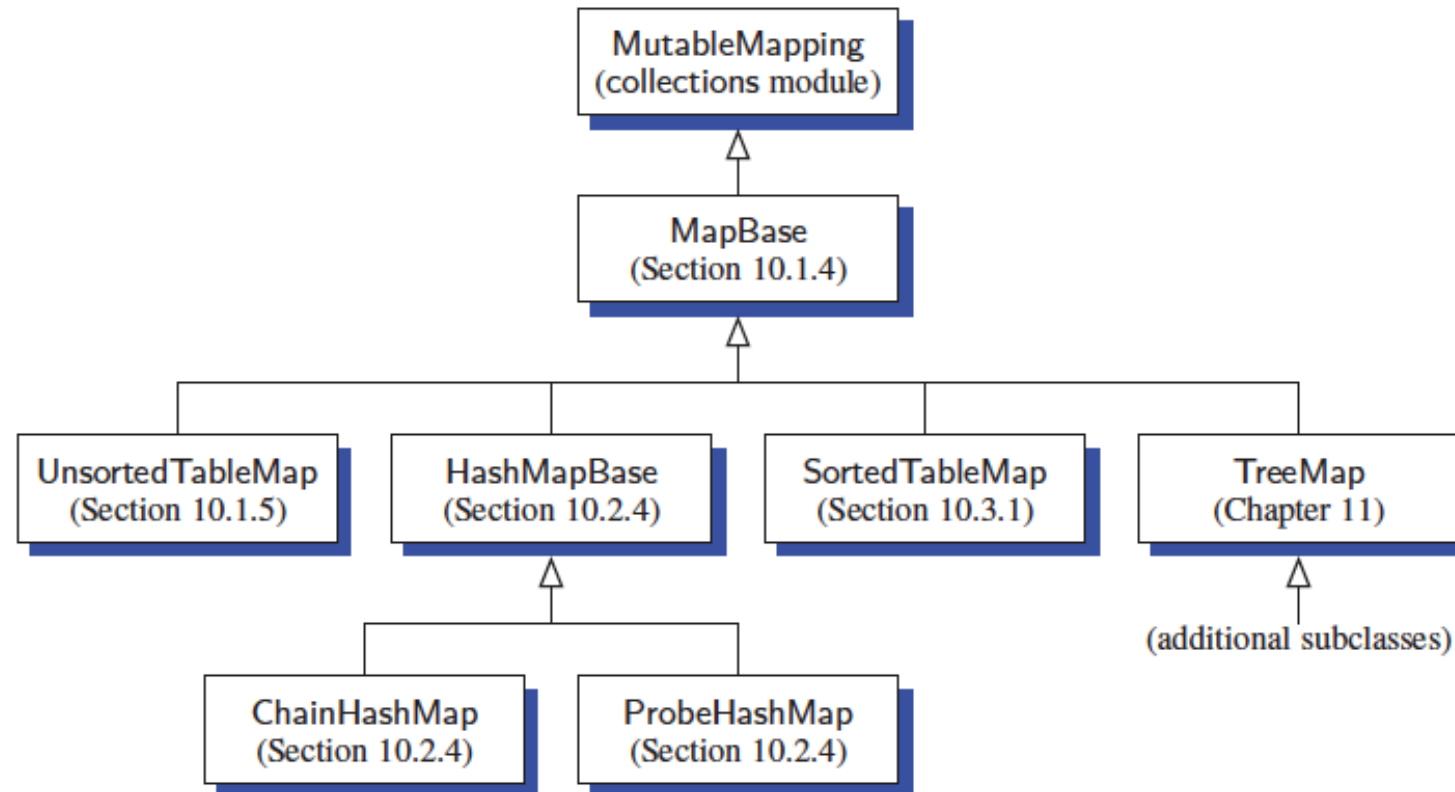


# A Simple List-Based Map

- 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



# Our MapBase Class



# The MapBase Abstract Class

```
1 class MapBase(MutableMapping):
2     """Our own abstract base class that includes a nonpublic _Item class."""
3
4     #----- nested _Item class -----
5     class _Item:
6         """Lightweight composite to store key-value pairs as map items."""
7         __slots__ = '_key', '_value'
8
9         def __init__(self, k, v):
10             self._key = k
11             self._value = v
12
13         def __eq__(self, other):
14             return self._key == other._key    # compare items based on their keys
15
16         def __ne__(self, other):
17             return not (self == other)        # opposite of __eq__
18
19         def __lt__(self, other):
20             return self._key < other._key     # compare items based on their keys
```

## Mapping(Collection)

<code>__getitem__</code> , <code>__iter__</code> , <code>__len__</code>	<code>__contains__</code> , <code>keys</code> , <code>items</code> , <code>values</code> , <code>get</code> , <code>__eq__</code> , and <code>__ne__</code>
---	--

## MutableMapping(Mapping)

<code>__getitem__</code> , <code>__setitem__</code> , <code>__delitem__</code> , <code>__iter__</code> , <code>__len__</code>	Inherited <code>Mapping</code> methods and <code>pop</code> , <code>popitem</code> , <code>clear</code> , <code>update</code> , and <code>setdefault</code>
---	---

# An Unsorted List Implementation

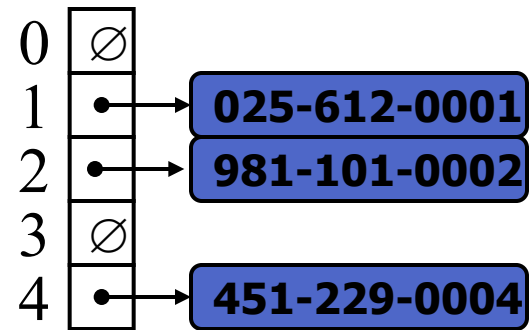
```
1 class UnsortedTableMap(MapBase):
2     """Map implementation using an unordered list."""
3
4     def __init__(self):
5         """Create an empty map."""
6         self._table = [ ]                # list of _Item's
7
8     def __getitem__(self, k):
9         """Return value associated with key k (raise KeyError if not found)."""
10        for item in self._table:
11            if k == item._key:
12                return item._value
13        raise KeyError('Key Error: ' + repr(k))
14
15    def __setitem__(self, k, v):
16        """Assign value v to key k, overwriting existing value if present."""
17        for item in self._table:
18            if k == item._key:                # Found a match:
19                item._value = v              # reassign value
20            return                          # and quit
21        # did not find match for key
22        self._table.append(self._Item(k,v))
23
```

```
24    def __delitem__(self, k):
25        """Remove item associated with key k (raise KeyError if not found)."""
26        for j in range(len(self._table)):
27            if k == self._table[j]._key:      # Found a match:
28                self._table.pop(j)          # remove item
29            return                          # and quit
30        raise KeyError('Key Error: ' + repr(k))
31
32    def __len__(self):
33        """Return number of items in the map."""
34        return len(self._table)
35
36    def __iter__(self):
37        """Generate iteration of the map's keys."""
38        for item in self._table:
39            yield item._key                  # yield the KEY
```

# Performance of a List-Based Map

- The unsorted list implementation is simple but not efficient.
  - `__getitem__`, `__setitem__`, and `__delitem__`, relies on a for loop to scan the underlying list of items in search of a matching key.
- Performance:
  - Inserting an item takes  $O(n)$  time since we first search for the existing keys and then insert the new item.
  - Searching for or removing an item takes  $O(n)$  time, since in the worst case (the item is not found) we traverse the entire list to look for an item with the given key.

# Hash Tables





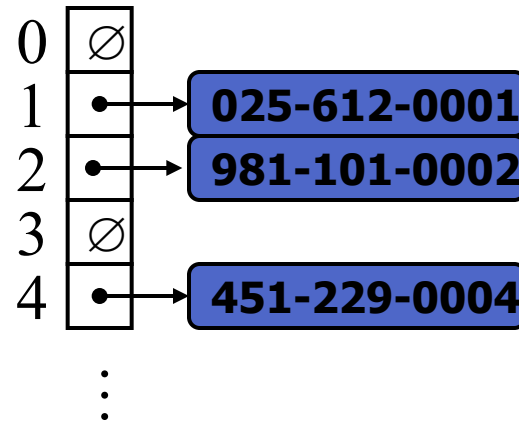
# Recall the notion of a Map

- Intuitively, a map  $M$  supports the abstraction of using keys as indices with a syntax such as  $M[k]$ .
- As a mental warm-up, consider a restricted setting in which a map with  $n$  items uses keys that are known to be integers in a range from 0 to  $N - 1$ , for some  $N \geq n$ .

0	1	2	3	4	5	6	7	8	9	10
	D		Z			C	Q			

# More General Kinds of Keys

- But what should we do if our keys are not integers in the range from 0 to  $N - 1$ ?
  - Use a **hash function** to map general keys to corresponding indices in a table.
  - For instance, the last four digits of a Social Security number.





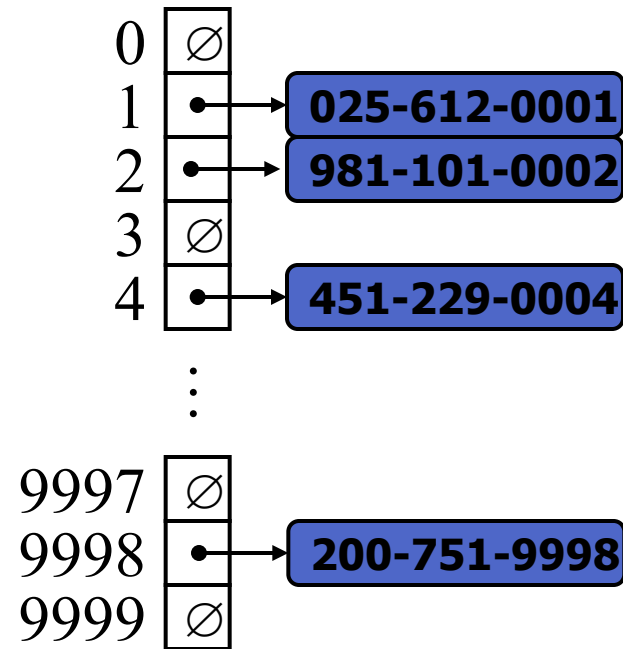
# Hash Functions and Hash Tables



- 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  $i = h(k)$

# SSN 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  
 $h(x) = \text{last four digits of } x$



# Hash Functions



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

Hash code:

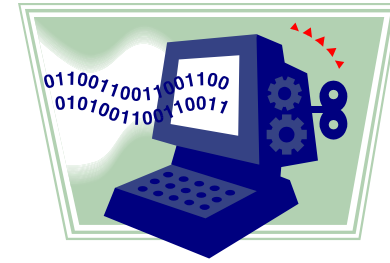
$h_1: \text{keys} \rightarrow \text{integers}$

Compression function:

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

- 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

# Hash Codes



- **Memory address:**
  - We reinterpret the memory address of the key object as an integer.
    - 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
- **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

# 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  $p(z) = p_{n-1}(z)$

# Compression Functions

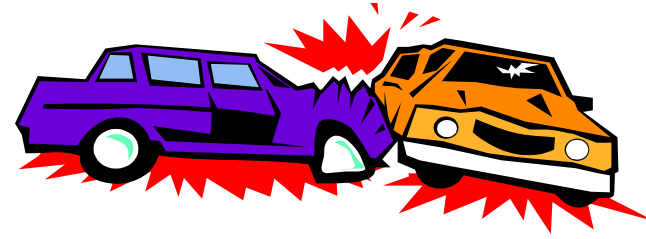


- Division:
  - $h_2(y) = y \bmod N$
  - The size  $N$  of the hash table is usually chosen to be a prime
  - The reason has to do with number theory and is beyond the scope of this course
- Multiply, Add and Divide (MAD):
  - $h_2(y) = [(ay + b) \bmod p] \bmod N$
  - $p$  is a prime number larger than  $N$
  - $a$  and  $b$  are nonnegative integers such that
$$a \bmod N \neq 0$$
$$a < p$$
  - This method distributes the keys in more in uniformly random way.

# Abstract Hash Table Class

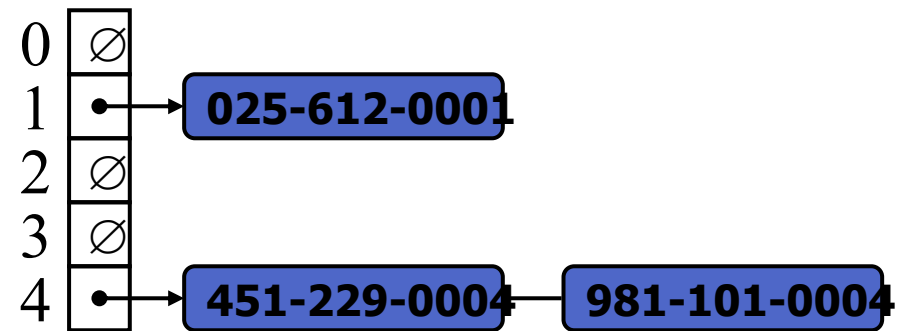
```
1 class HashMapBase(MapBase):
2     """Abstract base class for map using hash-table with MAD compression."""
3
4     def __init__(self, cap=11, p=109345121):
5         """Create an empty hash-table map."""
6         self._table = cap * [None]
7         self._n = 0                # number of entries in the map
8         self._prime = p            # prime for MAD compression
9         self._scale = 1 + randrange(p-1) # scale from 1 to p-1 for MAD
10        self._shift = randrange(p)    # shift from 0 to p-1 for MAD
11
12    def _hash_function(self, k):
13        return (hash(k)*self._scale + self._shift) % self._prime % len(self._table)
14
15    def __len__(self):
16        return self._n
17
18    def __getitem__(self, k):
19        j = self._hash_function(k)
20        return self._bucket_getitem(j, k)    # may raise KeyError
21
22    def __setitem__(self, k, v):
23        j = self._hash_function(k)
24        self._bucket_setitem(j, k, v)
25        if self._n > len(self._table) // 2:
26            self._resize(2 * len(self._table) - 1)    # number 2^x - 1 is often prime
27
28    def __delitem__(self, k):
29        j = self._hash_function(k)
30        self._bucket_delitem(j, k)    # may raise KeyError
31        self._n -= 1
32
33    def _resize(self, c):
34        old = list(self.items())
35        self._table = c * [None]
36        self._n = 0
37        for (k,v) in old:
38            self[k] = v    # reinsert old key-value pair
```

# Collision Handling



- 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 outside the table





# Map with Separate Chaining

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

**Algorithm** `get(k)`:  
**return** `A[h(k)].get(k)`

**Algorithm** `put(k,v)`:  
`t = A[h(k)].put(k,v)`  
**if** `t = null` **then** {k is a new key}  
    `n = n + 1`  
**return** `t`

**Algorithm** `remove(k)`:  
`t = A[h(k)].remove(k)`  
**if** `t ≠ null` **then** {k was found}  
    `n = n - 1`  
**return** `t`

# Hash Table with Separate Chaining

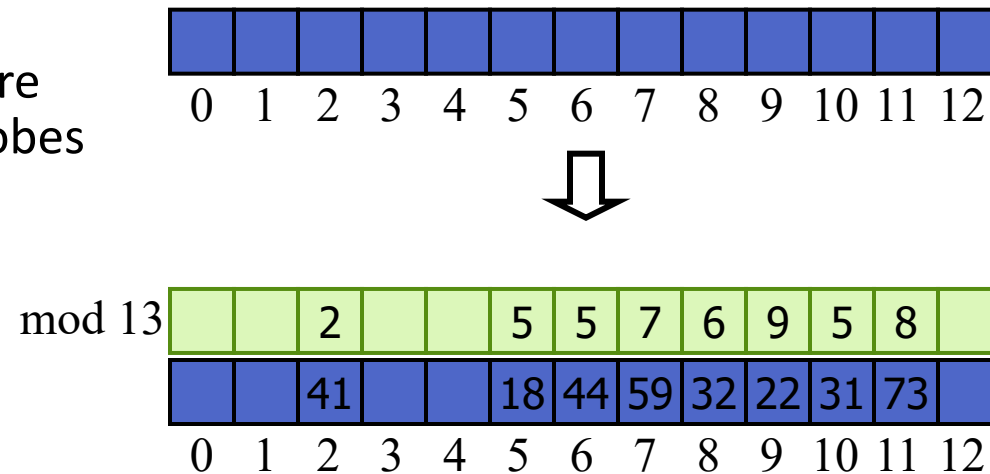
```
1 class ChainHashMap(HashMapBase):
2     """Hash map implemented with separate chaining for collision resolution."""
3
4     def _bucket_getitem(self, j, k):
5         bucket = self._table[j]
6         if bucket is None:
7             raise KeyError('Key Error: ' + repr(k))      # no match found
8         return bucket[k]                                  # may raise KeyError
9
10    def _bucket_setitem(self, j, k, v):
11        if self._table[j] is None:
12            self._table[j] = UnsortedTableMap( )          # bucket is new to the table
13            oldsize = len(self._table[j])
14            self._table[j][k] = v
15            if len(self._table[j]) > oldsize:               # key was new to the table
16                self._n += 1                                # increase overall map size
17
18    def _bucket_delitem(self, j, k):
19        bucket = self._table[j]
20        if bucket is None:
21            raise KeyError('Key Error: ' + repr(k))      # no match found
22        del bucket[k]                                       # may raise KeyError
23
24    def __iter__(self):
25        for bucket in self._table:
26            if bucket is not None:                          # a nonempty slot
27                for key in bucket:
28                    yield key
```

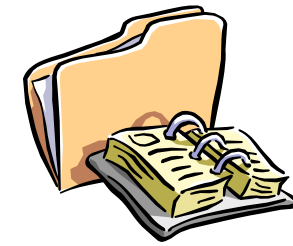
# Linear Probing

- 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





# Search with Linear Probing

- Consider a hash table  $A$  that uses linear probing
- $\text{get}(k)$ 
  - We start at cell  $h(k)$
  - We probe consecutive locations until one of the following occurs
    - An item with key  $k$  is found, or
    - An empty cell is found, or
    - $N$  cells have been unsuccessfully probed

## Algorithm $\text{get}(k)$

```
 $i \leftarrow h(k)$   
 $p \leftarrow 0$   
repeat  
     $c \leftarrow A[i]$   
    if  $c = \emptyset$   
        return null  
    else if  $c.\text{getKey}() = k$   
        return  $c.\text{getValue}()$   
    else  
         $i \leftarrow (i + 1) \bmod N$   
         $p \leftarrow p + 1$   
until  $p = N$   
return null
```

# Updates with Linear Probing

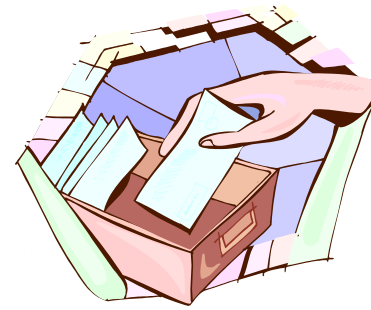
- To handle insertions and deletions, we introduce a special object, called *AVAILABLE*, which replaces deleted elements
- **remove( $k$ )**
  - We search for an entry with key  $k$
  - If such an entry  $(k, o)$  is found, we replace it with the special item *AVAILABLE* and we return element  $o$
  - Else, we return *null*
- **put( $k, o$ )**
  - We throw an exception if the table is full
  - We start at cell  $h(k)$
  - We probe consecutive cells until one of the following occurs
    - A cell  $i$  is found that is either empty or stores *AVAILABLE*, or
    - $N$  cells have been unsuccessfully probed
  - We store  $(k, o)$  in cell  $i$

# Hash Table with Linear Probing

```
1 class ProbeHashMap(HashMapBase):
2     """Hash map implemented with linear probing for collision resolution."""
3     _AVAIL = object() # sentinel marks locations of previous deletions
4
5     def _is_available(self, j):
6         """Return True if index j is available in table."""
7         return self._table[j] is None or self._table[j] is ProbeHashMap._AVAIL
8
9     def _find_slot(self, j, k):
10        """Search for key k in bucket at index j.
11
12        Return (success, index) tuple, described as follows:
13        If match was found, success is True and index denotes its location.
14        If no match found, success is False and index denotes first available slot.
15        """
16        firstAvail = None
17        while True:
18            if self._is_available(j):
19                if firstAvail is None:
20                    firstAvail = j # mark this as first avail
21                if self._table[j] is None:
22                    return (False, firstAvail) # search has failed
23                elif k == self._table[j]._key:
24                    return (True, j) # found a match
25                j = (j + 1) % len(self._table) # keep looking (cyclically)
```

```
26     def _bucket_getitem(self, j, k):
27         found, s = self._find_slot(j, k)
28         if not found:
29             raise KeyError('Key Error: ' + repr(k)) # no match found
30         return self._table[s]._value
31
32     def _bucket_setitem(self, j, k, v):
33         found, s = self._find_slot(j, k)
34         if not found:
35             self._table[s] = self._Item(k,v) # insert new item
36             self._n += 1 # size has increased
37         else:
38             self._table[s]._value = v # overwrite existing
39
40     def _bucket_delitem(self, j, k):
41         found, s = self._find_slot(j, k)
42         if not found:
43             raise KeyError('Key Error: ' + repr(k)) # no match found
44             self._table[s] = ProbeHashMap._AVAIL # mark as vacated
45
46     def __iter__(self):
47         for j in range(len(self._table)): # scan entire table
48             if not self._is_available(j):
49                 yield self._table[j]._key
```

# Double Hashing



- Double hashing uses a secondary hash function  $d(k)$  and handles collisions by placing an item in the first available cell of the series

$$(i + jd(k)) \bmod N$$

for  $j = 0, 1, \dots, N-1$

- The secondary hash function  $d(k)$  cannot have zero values
- The table size  $N$  must be a prime to allow probing of all the cells

- Common choice of compression function for the secondary hash function:

$$d_2(k) = q - k \bmod q$$

where

- $q < N$
- $q$  is a prime
- The possible values for  $d_2(k)$  are  $1, 2, \dots, q$

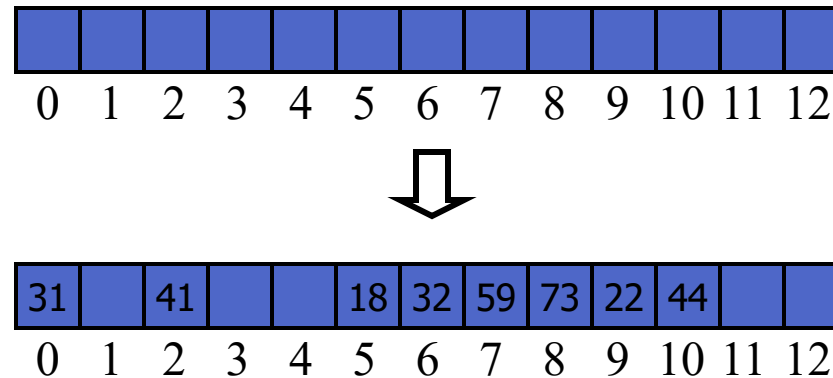
# Example of Double Hashing

- Consider a hash table storing integer keys that handles collision with double hashing

- $N = 13$
- $h(k) = k \bmod 13$
- $d(k) = 7 - k \bmod 7$

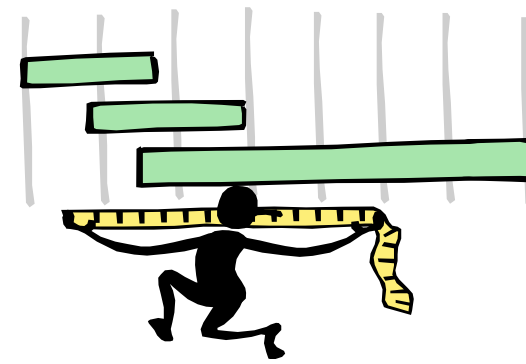
- Insert keys 18, 41, 22, 44, 59, 32, 31, 73, in this order

$k$	$h(k)$	$d(k)$	Probes	
18	5	3	5	
41	2	1	2	
22	9	6	9	
44	5	5	5	10
59	7	4	7	
32	6	3	6	
31	5	4	5	9 0
73	8	4	8	





# Performance of Hashing



- In the worst case, searches, insertions and removals on a hash table take  $O(n)$  time
- The worst case occurs when all the keys inserted into the map collide
- The load factor  $\alpha = n/N$  affects the performance of a hash table
- Assuming that the hash values are like random numbers, it can be shown that the expected number of probes for an insertion with open addressing is
$$1 / (1 - \alpha)$$
- The expected running time of all the dictionary ADT operations in a hash table is  $O(1)$
- In practice, hashing is very fast provided the load factor is not close to 100%
- Applications of hash tables:
  - small databases
  - compilers
  - browser caches