

ROBERT SEDGEWICK | KEVIN WAYNE

<https://algs4.cs.princeton.edu>

3.4 HASH TABLES

- ▶ *hash functions*
- ▶ *separate chaining*
- ▶ *linear probing*
- ▶ *context*

Symbol table implementations: summary

implementation	guarantee			average case			ordered ops?	key interface
	search	insert	delete	search	insert	delete		
sequential search (unordered list)	n	n	n	n	n	n		equals()
binary search (ordered array)	$\log n$	n	n	$\log n$	n	n	✓	compareTo()
BST	n	n	n	$\log n$	$\log n$	\sqrt{n}	✓	compareTo()
red-black BST	$\log n$	$\log n$	$\log n$	$\log n$	$\log n$	$\log n$	✓	compareTo()
hashing	n	n	n	1^\dagger	1^\dagger	1^\dagger		equals() hashCode()

† under suitable technical assumptions

Q. Can we do better?

A. Yes, but only with different access to the data.

Hashing: basic plan

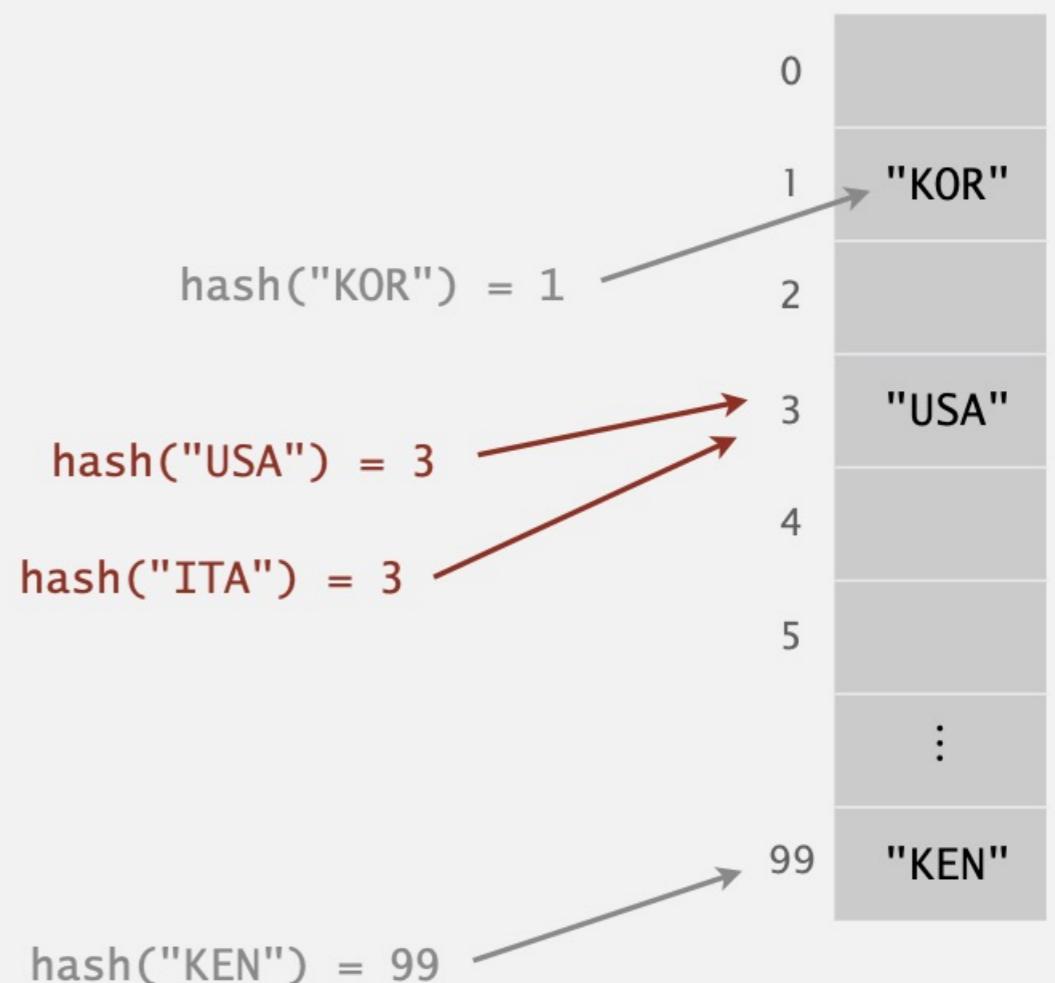
Save key-value pairs in a **key-indexed table** (index is a function of the key).

Hash function. Function that maps a key to an array index.

Collision. Two distinct keys that hash to same index.

Issue. Collisions are inevitable.

- How to limit collisions?
- How to accommodate collisions?





Algorithms

ROBERT SEDGEWICK | KEVIN WAYNE

<https://algs4.cs.princeton.edu>

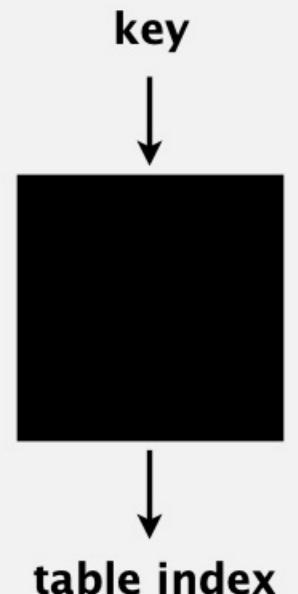
3.4 HASH TABLES

- ▶ ***hash functions***
- ▶ ***separate chaining***
- ▶ ***linear probing***
- ▶ ***context***

Designing a hash function

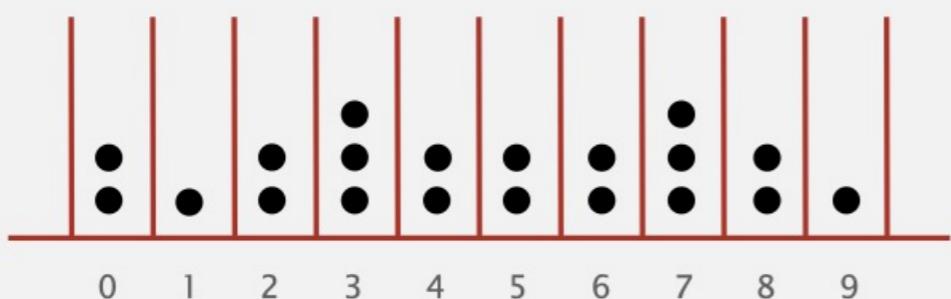
Required properties. [for correctness]

- Deterministic.
- Each key hashes to a table index between 0 and $m - 1$.

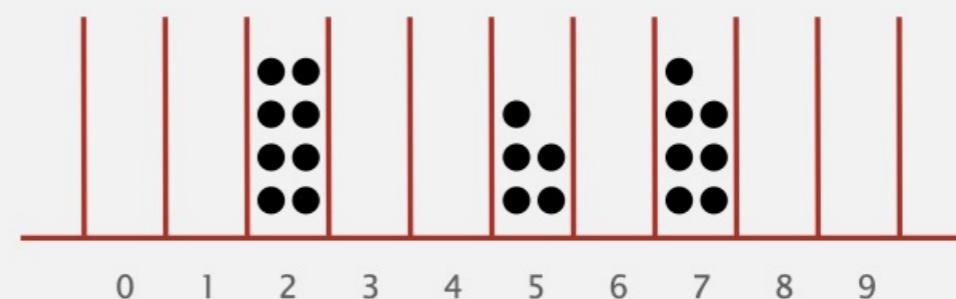


Desirable properties. [for performance]

- Very fast to compute. ← constants matter
- For any subset of n input keys, each table index gets approximately n / m keys.



leads to good hash-table performance
($m = 10, n = 20$)

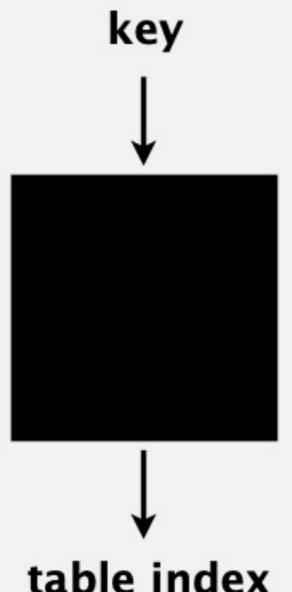


leads to bad hash-table performance
($m = 10, n = 20$)

Designing a hash function

Required properties. [for correctness]

- Deterministic.
- Each key hashes to a table index between 0 and $m - 1$.

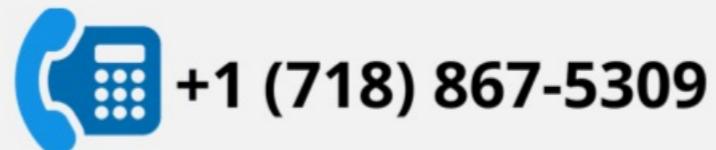
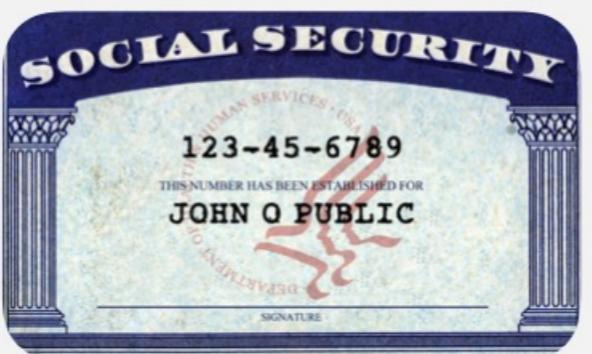


Desirable properties. [for performance]

- Very fast to compute.
- For any subset of n input keys, each table index gets approximately n / m keys.

Ex 1. Last 4 digits of U.S. Social Security number.

Ex 2. Last 4 digits of phone number.

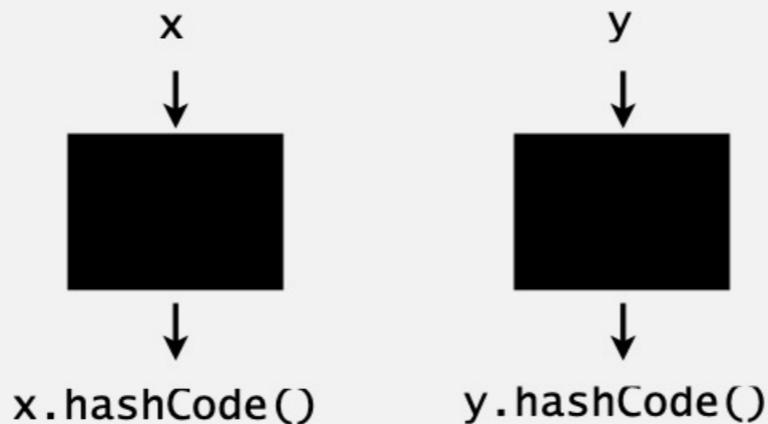


Java's hashCode() conventions

All Java classes inherit a method `hashCode()`, which returns a 32-bit `int`.

Requirement. If `x.equals(y)`, then `(x.hashCode() == y.hashCode())`.

Highly desirable. If `!x.equals(y)`, then `(x.hashCode() != y.hashCode())`.



Customized implementations. `Integer`, `Double`, `String`, `java.net.URL`, ...

Legal (but undesirable) implementation. Always return 17.

User-defined types. Users are on their own.

Implementing hashCode(): integers and doubles

Java library implementations

```
public final class Integer
{
    private final int value;
    ...
    public int hashCode()
    { return value; }
}
```

```
public final class Double
{
    private final double value;
    ...
    public int hashCode()
    {
        long bits = doubleToLongBits(value);
        return (int) (bits ^ (bits >>> 32));
    }
}
```

convert to IEEE 64-bit representation;
xor most significant 32-bits
with least significant 32-bits

Implementing hashCode(): arrays

[31x + y rule.](#)

- Initialize hash to 1.
- Repeatedly multiply hash by 31 and add next integer in array.

```
public class Arrays
{
    ...

    public static int hashCode(int[] a) {
        if (a == null)
            return 0; ← special case for null

        int hash = 1;
        for (int i = 0; i < a.length; i++)
            hash = 31*hash + a[i];
        return hash;
    }
}
```

Java library implementation

prime
↓

31x + y rule

Implementing hashCode(): user-defined types

```
public final class Transaction
{
    private final String who;
    private final Date when;
    private final double amount;

    public Transaction(String who, Date when, double amount)
    { /* as before */ }

    public boolean equals(Object y)
    { /* as before */ }

    ...

    public int hashCode()
    {
        int hash = 1;
        hash = 31*hash + who.hashCode();
        hash = 31*hash + when.hashCode();
        hash = 31*hash + ((Double) amount).hashCode();
        return hash;
    }
}
```

for reference types,
use hashCode()
for primitive types,
use hashCode()
of wrapper type

Implementing hashCode(): user-defined types

```
public final class Transaction
{
    private final String who;
    private final Date when;
    private final double amount;

    public Transaction(String who, Date when, double amount)
    { /* as before */ }

    public boolean equals(Object y)
    { /* as before */ }

    ...

    public int hashCode()
    {
        return Objects.hash(who, when, amount); ← shorthand
    }
}
```

Implementing hashCode()

“Standard” recipe for user-defined types.

- Combine each significant field using the $31x + y$ rule.
- Shortcut 1: use `Objects.hash()` for all fields (except arrays).
- Shortcut 2: use `Arrays.hashCode()` for primitive arrays.
- Shortcut 3: use `Arrays.deepHashCode()` for object arrays.



Principle. Every significant field contributes to hash.

In practice. Recipe above works reasonably well; used in Java libraries.



Which function maps hashable keys to integers between 0 and $m-1$?

A.

```
private int hash(Key key)
{   return key.hashCode() % m; }
```

B.

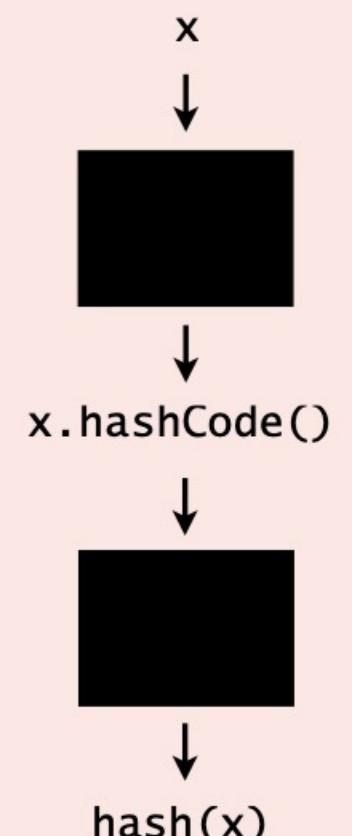
```
private int hash(Key key)
{   return Math.abs(key.hashCode()) % m; }
```

C.

Both A and B.

D.

Neither A nor B.



Modular hashing

Hash code. An int between -2^{31} and $2^{31} - 1$.

Hash function. An int between 0 and $m - 1$ (for use as array index).

typically a prime or power of 2

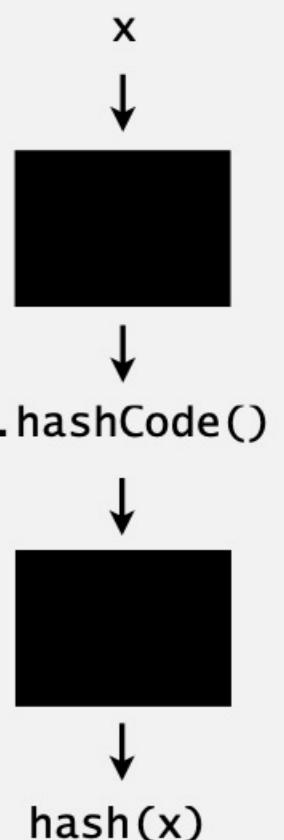
```
private int hash(Key key)
{   return key.hashCode() % m; }
```

bug

```
private int hash(Key key)
{   return Math.abs(key.hashCode()) % m; }
```

1-in-a-billion bug

hashCode() of "polygenelubricants" is -2^{31}



```
private int hash(Key key)
{   return (key.hashCode() & 0xffffffff) % m; }
```

correct

if m is a power of 2, can use
 $\text{key.hashCode()} \& (m-1)$

Uniform hashing assumption

Uniform hashing assumption. Any key is equally likely to hash to one of m possible indices.

and independently of other keys

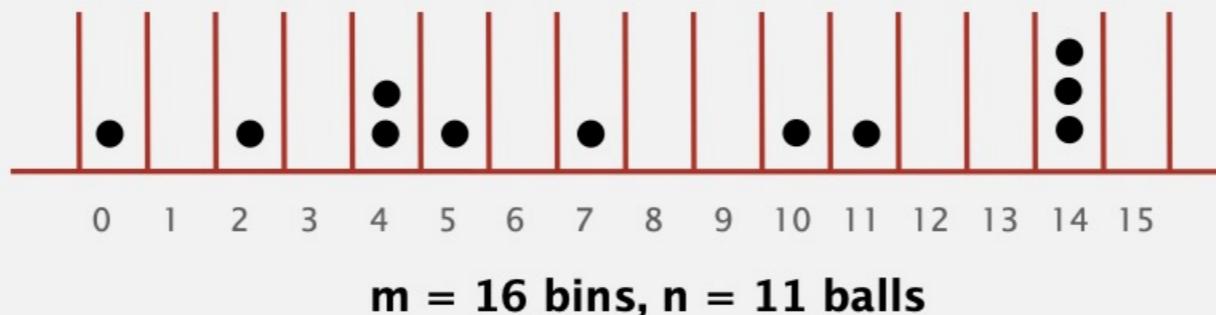


Uniform hashing assumption

Uniform hashing assumption. Any key is equally likely to hash to one of m possible indices.

and independently of other keys

Bins and balls. Toss n balls uniformly at random into m bins.



Bad news. [birthday problem]

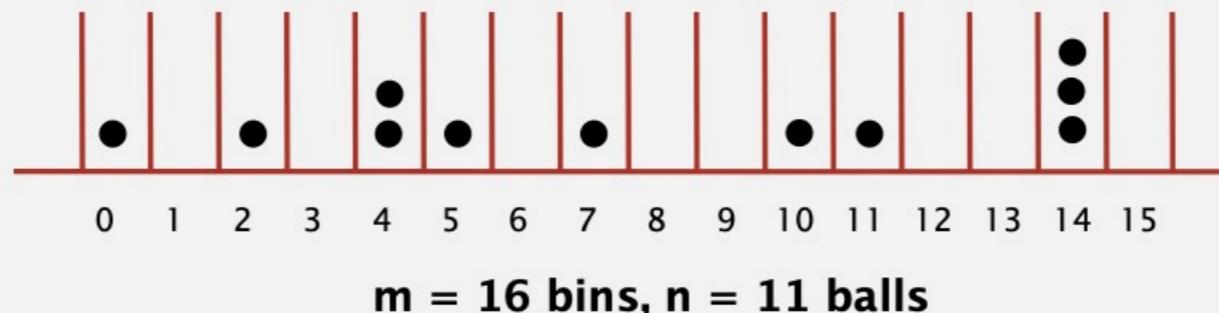
- In a random group of 23 people, more likely than not that two people share the same birthday.
- Expect two balls in the same bin after $\sim \sqrt{\pi m / 2}$ tosses.

23.9 when $m = 365$

Uniform hashing assumption

Uniform hashing assumption. Any key is equally likely to hash to one of m possible indices.

Bins and balls. Toss n balls uniformly at random into m bins.



Good news. [load balancing]

- When $n \gg m$, expect most bins to have approximately n / m balls.
- When $n = m$, expect most loaded bin has $\sim \ln n / \ln \ln n$ balls.

Binomial($n, 1 / m$)



Algorithms

ROBERT SEDGEWICK | KEVIN WAYNE

<https://algs4.cs.princeton.edu>

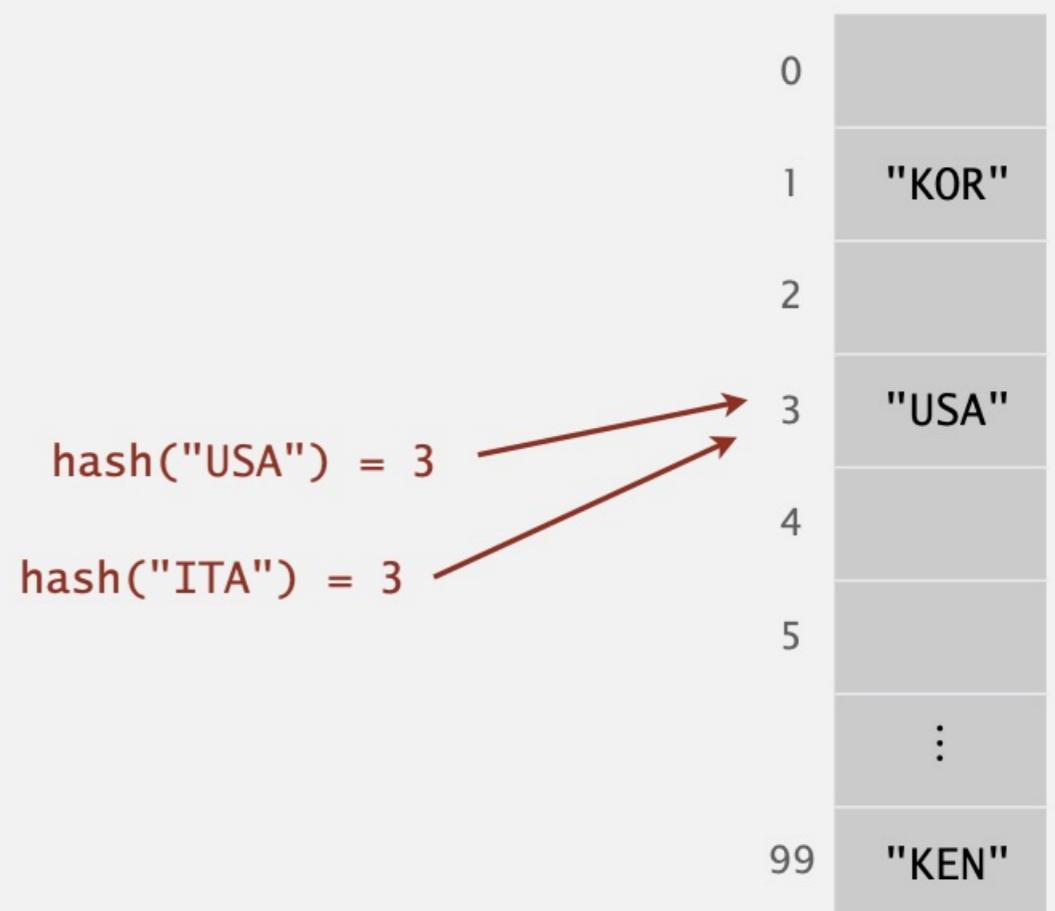
3.4 HASH TABLES

- ▶ *hash functions*
- ▶ ***separate chaining***
- ▶ *linear probing*
- ▶ *context*

Collisions

Collision. Two distinct keys that hash to the same index.

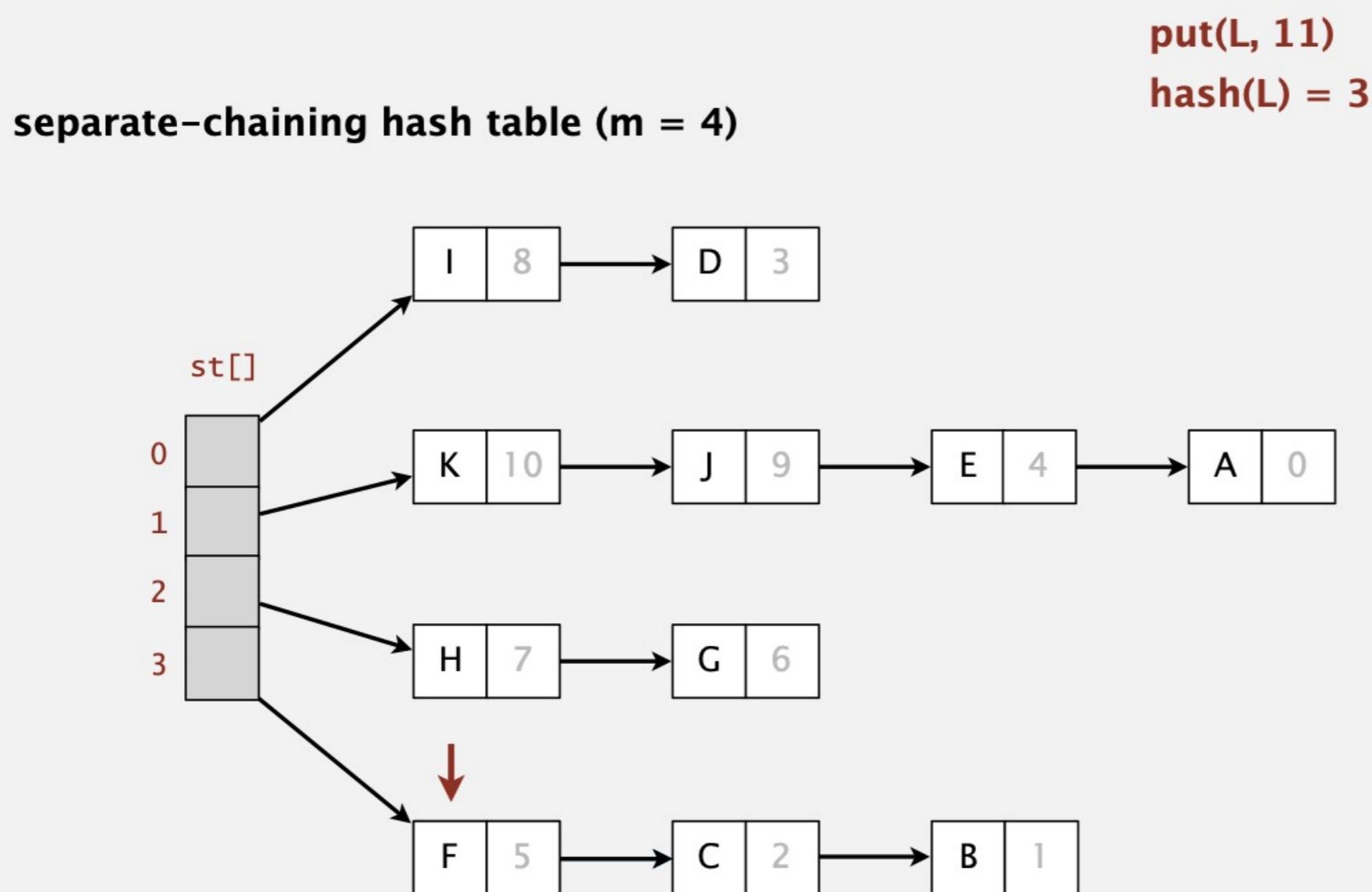
- Birthday problem \Rightarrow can't avoid collisions. ← unless you have a ridiculous (quadratic) amount of memory
- Load balancing \Rightarrow no index gets too many collisions.
 \Rightarrow ok to scan through all colliding keys.



Separate-chaining hash table

Use an array of m linked lists.

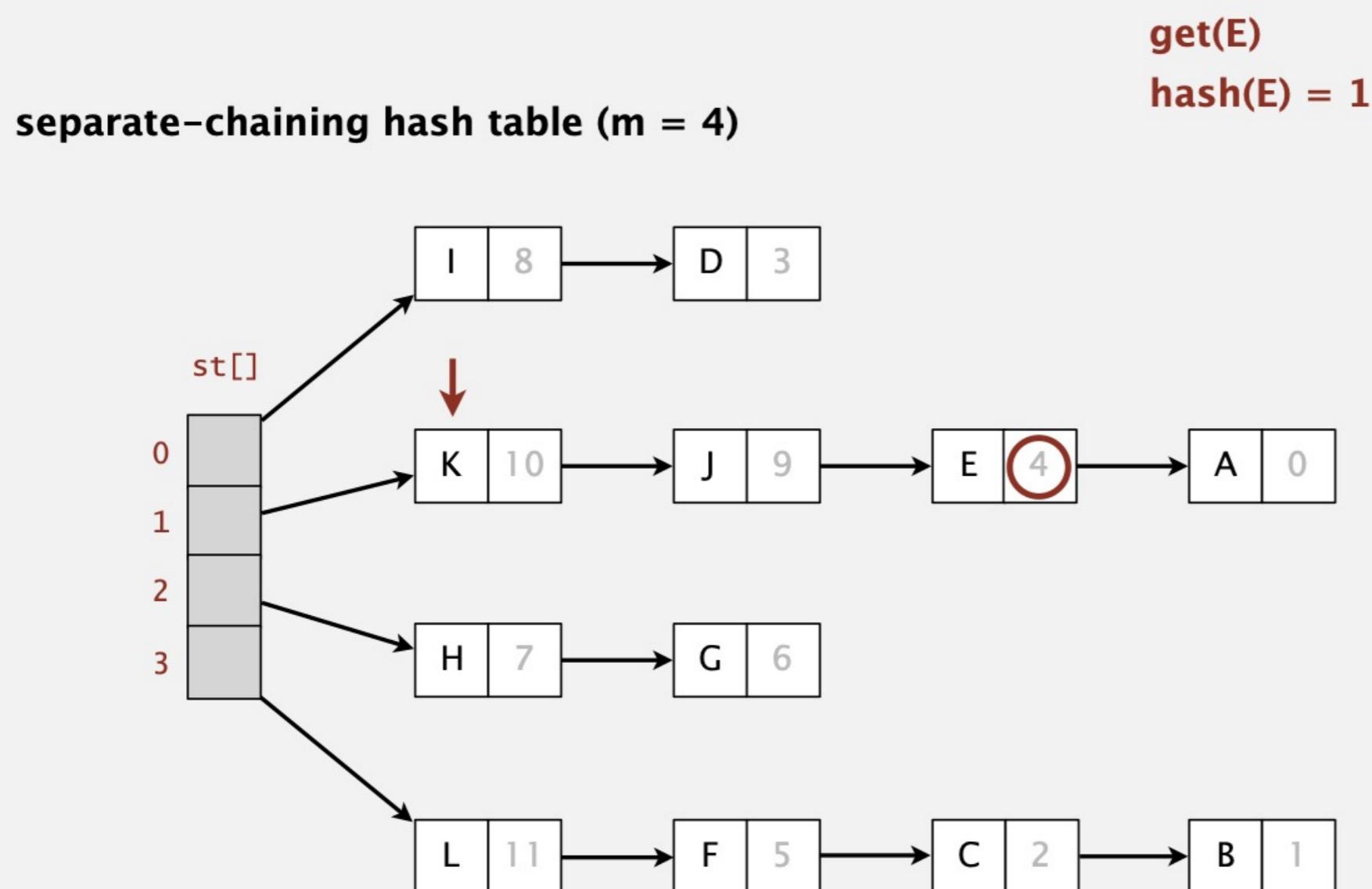
- Hash: map key to table index i between 0 and $m - 1$.
- Insert: add key-value pair at front of chain i (if not already in chain).



Separate-chaining hash table

Use an array of m linked lists.

- Hash: map key to table index i between 0 and $m - 1$.
- Insert: add key-value pair at front of chain i (if not already in chain).
Search: perform sequential search in chain i .



Separate-chaining hash table: Java implementation

```
public class SeparateChainingHashST<Key, Value>
{
    private int m = 128;                      // number of chains
    private Node[] st = new Node[m];           // array of chains
                                                ← array resizing
                                                code omitted

    private static class Node
    {
        private Object key;                   ← no generic array creation
        private Object val;                 (declare key and value of type Object)
        private Node next;
        ...
    }

    private int hash(Key key)
    {   return (key.hashCode() & 0xffffffff) % m;   }

    public Value get(Key key) {
        int i = hash(key);
        for (Node x = st[i]; x != null; x = x.next)
            if (key.equals(x.key)) return (Value) x.val;
        return null;
    }
}
```

Separate-chaining hash table: Java implementation

```
public class SeparateChainingHashST<Key, Value>
{
    private int m = 128;                      // number of chains
    private Node[] st = new Node[m]; // array of chains

    private static class Node
    {
        private Object key;
        private Object val;
        private Node next;
        ...
    }

    private int hash(Key key)
    {   return (key.hashCode() & 0xffffffff) % m;   }

    public void put(Key key, Value val)
    {
        int i = hash(key);
        for (Node x = st[i]; x != null; x = x.next)
            if (key.equals(x.key)) { x.val = val; return; }
        st[i] = new Node(key, val, st[i]);
    }

}
```

Analysis of separate chaining

Recall load balancing. Under uniform hashing assumption, length of each chain is tightly concentrated around mean = n / m .



Consequence. Expected number of **probes** for search/insert is $\Theta(n / m)$.

- m too small \Rightarrow chains too long.
- m too large \Rightarrow too many empty chains.
- Typical choice: $m \sim \frac{1}{4}n \Rightarrow \Theta(1)$ time for search/insert.

calls to either
equals() or hashCode()



↑
 m times faster than
sequential search

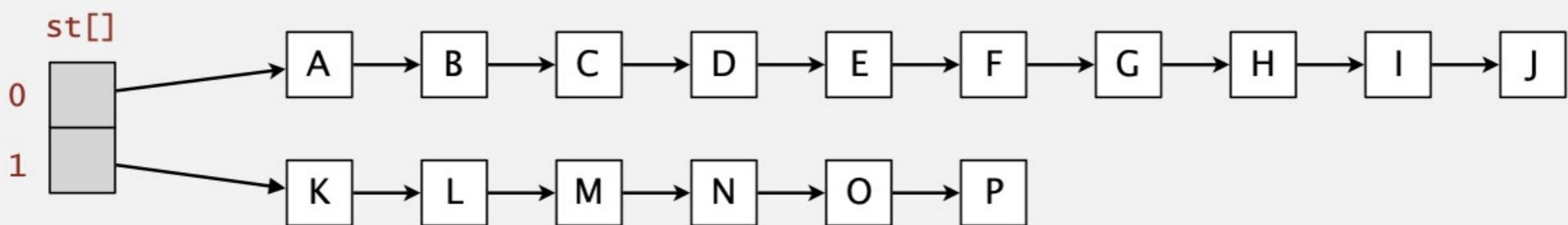
Resizing in a separate-chaining hash table

Goal. Average length of list $n / m = \text{constant}$.

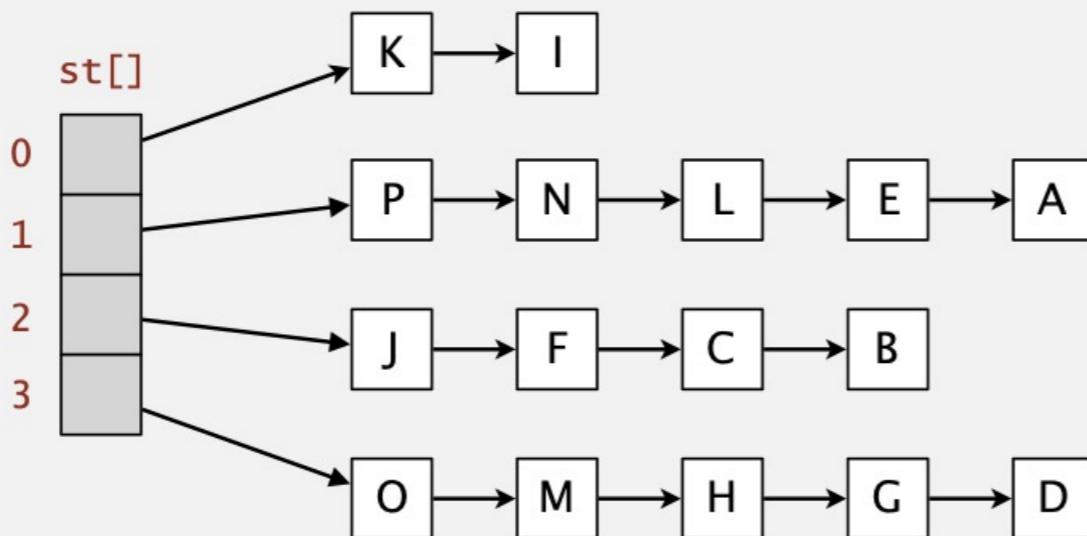
- Double length m of array when $n / m \geq 8$.
- Halve length m of array when $n / m \leq 2$.
- Note: need to rehash all keys when resizing.

x.hashCode() does not change;
but hash(x) typically does

before resizing ($n/m = 8$)



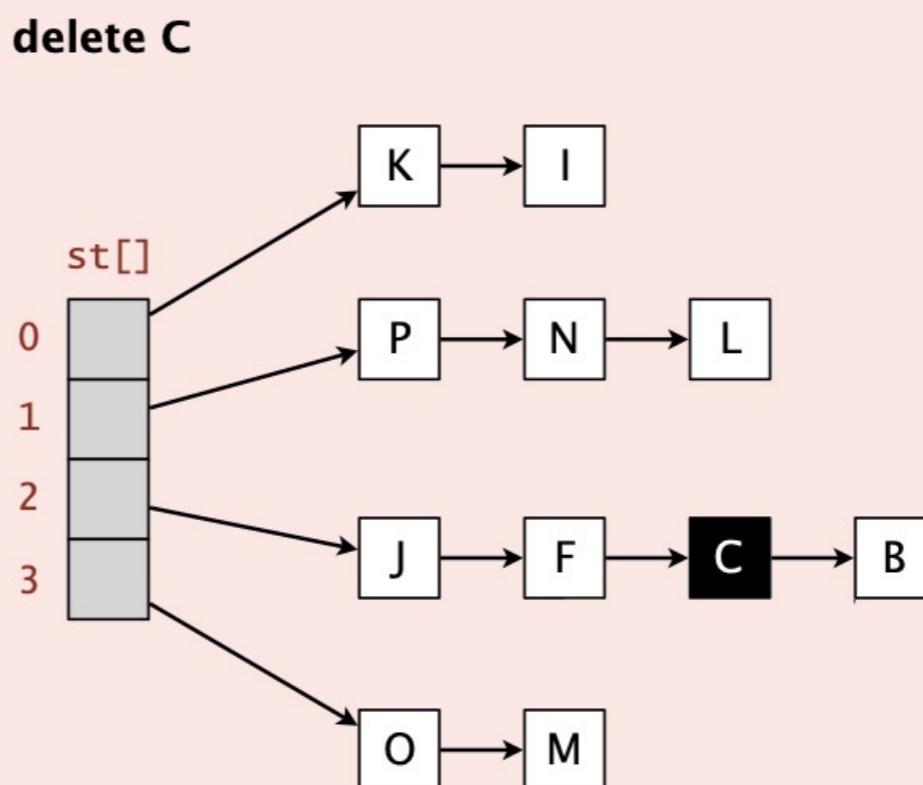
after resizing ($n/m = 4$)





How to delete a key-value pair from a separate-chaining hash table?

- A. Search for key; remove key-value pair from linked list.
- B. Compute hash of key; reinsert all other key-value pairs in chain.
- C. Either A or B.
- D. Neither A nor B.



Symbol table implementations: summary

implementation	guarantee			average case			ordered ops?	key interface
	search	insert	delete	search	insert	delete		
sequential search (unordered list)	n	n	n	n	n	n		equals()
binary search (ordered array)	$\log n$	n	n	$\log n$	n	n	✓	compareTo()
BST	n	n	n	$\log n$	$\log n$	\sqrt{n}	✓	compareTo()
red-black BST	$\log n$	$\log n$	$\log n$	$\log n$	$\log n$	$\log n$	✓	compareTo()
separate chaining	n	n	n	1^\dagger	1^\dagger	1^\dagger		equals() hashCode()

† under uniform hashing assumption

Algorithms

ROBERT SEDGEWICK | KEVIN WAYNE

<https://algs4.cs.princeton.edu>

3.4 HASH TABLES

- ▶ *hash functions*
- ▶ *separate chaining*
- ▶ ***linear probing***
- ▶ *context*

Linear-probing hash table

- Maintain key–value pairs in two parallel arrays, with one key per cell.
- Resolve collisions by probing: search successive cells until either finding the key or an unused cell.

Inserting into a linear-probing hash table.

linear-probing hash table	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
keys[]	P	M			A	C		H	L		E			R	X	
put(K, 14)								K								
hash(K) = 7									14							
vals[]	11	10			9	5		6	12		13			4	8	

Linear-probing hash table

- Maintain key-value pairs in two parallel arrays, with one key per cell.
- Resolve collisions by probing: search successive cells until either finding the key or an unused cell.

Searching in a linear-probing hash table.

linear-probing hash table

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
keys[]	P	M			A	C		H	L	K	E				R	X
	get(K)				get(Z)					K	Z					
	hash(K) = 7				hash(Z) = 8											
vals[]	11	10			9	5		6	12	14	13				4	8

Linear-probing hash table demo

Hash. Map key to integer i between 0 and $m - 1$.

Insert. Put at table index i if free; if not try $i + 1, i + 2, \dots$

Search. Search table index i ; if occupied but no match, try $i + 1, i + 2, \dots$

Note. Array length m **must** be greater than number of key-value pairs n .

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
keys[]	P	M			A	C	S	H	L		E				R	X
m =	16															



Linear-probing symbol table: Java implementation

```
public class LinearProbingHashST<Key, Value>
{
    private int m = 32768;
    private Value[] vals = (Value[]) new Object[m];
    private Key[] keys = (Key[]) new Object[m];

    private int hash(Key key)
    {   return (key.hashCode() & 0xffffffff) % m; }

    private void put(Key key, Value val) { /* next slide */ }

    public Value get(Key key)
    {
        for (int i = hash(key); keys[i] != null; i = (i+1) % m)
            if (key.equals(keys[i]))
                return vals[i];
        return null;
    }
}
```

array resizing
code omitted

Linear-probing symbol table: Java implementation

```
public class LinearProbingHashST<Key, Value>
{
    private int m = 32768;
    private Value[] vals = (Value[]) new Object[m];
    private Key[] keys = (Key[]) new Object[m];

    private int hash(Key key)
    { return (key.hashCode() & 0xffffffff) % m; }

    private Value get(Key key) { /* prev slide */ }

    public void put(Key key, Value val)
    {
        int i;
        for (i = hash(key); keys[i] != null; i = (i+1) % m)
            if (keys[i].equals(key))
                break;
        keys[i] = key;
        vals[i] = val;
    }
}
```



Under the uniform hashing assumption, where is the next key most likely to be added in this linear-probing hash table (no resizing)?

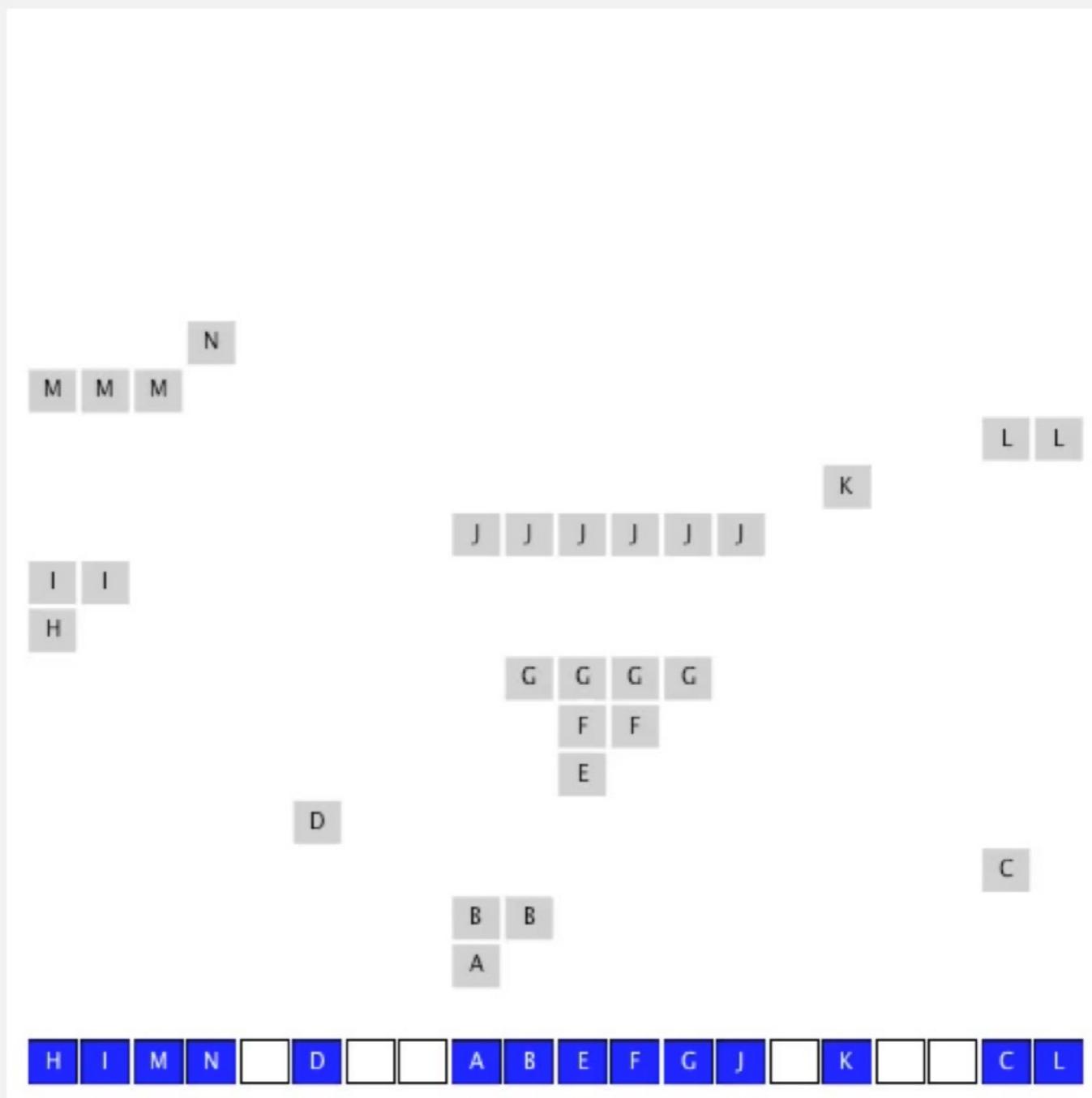
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
H	I	M	N		D			A	B	E	F	G	J		K			C	L

- A. Index 7.
- B. Index 14.
- C. Either index 4 or 14.
- D. All open indices are equally likely.

Clustering

Cluster. A contiguous block of keys.

Observation. New keys disproportionately likely to hash into big clusters.



Analysis of linear probing

Proposition. Under uniform hashing assumption, the average # of probes in a linear-probing hash table of size m that contains $n = \alpha m$ keys is at most

$$\frac{1}{2} \left(1 + \frac{1}{1 - \alpha} \right)$$

search hit

$$\frac{1}{2} \left(1 + \frac{1}{(1 - \alpha)^2} \right)$$

search miss / insert

Pf. [beyond course scope]

[My first analysis of an algorithm, originally done during summer 1962 in Madison.] 8354

NOTES ON "OPEN" ADDRESSING. D. Knuth, 7/22/63

1. Introduction and Definitions. Open addressing is a widely-used technique for keeping "symbol tables." The method was first used in 1954 by Samuel, Amdahl, and Boehme in an assembly program for the IBM 701. An extensive discussion of the method was given by Peterson in 1957 [1], and frequent references have been made to it ever since (e.g. Schay and Spruth [2], Iverson [3]). However, the timing characteristics have apparently never been exactly established, and indeed the author has heard reports of several reputable mathematicians who failed to find the solution after some trial. Therefore it is the purpose of this note to indicate one way by which the solution can be obtained.



Parameters.

- m too large \Rightarrow too many empty array entries.
- m too small \Rightarrow search time blows up.
- Typical choice: $\alpha = n/m \sim 1/2.$

probes for search hit is about 3/2
probes for search miss is about 5/2

Resizing in a linear-probing hash table

Goal. Average length of list $n / m \leq \frac{1}{2}$.

- Double length of array m when $n / m \geq \frac{1}{2}$.
- Halve length of array m when $n / m \leq \frac{1}{8}$.
- Need to rehash all keys when resizing.

before resizing

	0	1	2	3	4	5	6	7
keys[]		E	S		R	A		
vals[]		1	0			3	2	

after resizing

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
keys[]				A		S			E				R			
vals[]					2		0			1				3		



How to delete a key-value pair from a linear-probing hash table?

- A. Search for key; remove key-value pair from arrays.
- B. Search for key; remove key-value pair from arrays.
Shift all keys in **cluster** after deleted key 1 position to left.
- C. Either A and B.
- D. Neither A nor B.

		cluster after deleted key															
		0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
before deleting S	keys[]	P	M			A	C	S	H	L		E			R	X	
	vals[]	10	9			8	4	0	5	11		12			3	7	

ST implementations: summary

implementation	guarantee			average case			ordered ops?	key interface
	search	insert	delete	search	insert	delete		
sequential search (unordered list)	n	n	n	n	n	n		equals()
binary search (ordered array)	$\log n$	n	n	$\log n$	n	n	✓	compareTo()
BST	n	n	n	$\log n$	$\log n$	\sqrt{n}	✓	compareTo()
red-black BST	$\log n$	$\log n$	$\log n$	$\log n$	$\log n$	$\log n$	✓	compareTo()
separate chaining	n	n	n	1^\dagger	1^\dagger	1^\dagger		equals() hashCode()
linear probing	n	n	n	1^\dagger	1^\dagger	1^\dagger		equals() hashCode()

† under uniform hashing assumption