

## 3.4 Hash Tables



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

# Optimize judiciously

*“ More computing sins are committed in the name of efficiency (without necessarily achieving it) than for any other single reason—including blind stupidity. ” — William A. Wulf*

*“ We should forget about small efficiencies, say about 97% of the time: premature optimization is the root of all evil. ” — Donald E. Knuth*

*“ We follow two rules in the matter of optimization:  
Rule 1: Don't do it.  
Rule 2 (for experts only). Don't do it yet - that is, not until you have a perfectly clear and unoptimized solution. ” — M. A. Jackson*

**Reference: Effective Java by Joshua Bloch**



# ST implementations: summary

implementation	guarantee			average case			ordered iteration?	operations on keys
	search	insert	delete	search hit	insert	delete		
sequential search (linked list)	N	N	N	N/2	N	N/2	no	<code>equals()</code>
binary search (ordered array)	$\lg N$	N	N	$\lg N$	N/2	N/2	yes	<code>compareTo()</code>
BST	N	N	N	$1.38 \lg N$	$1.38 \lg N$	?	yes	<code>compareTo()</code>
red-black tree	$2 \lg N$	$2 \lg N$	$2 \lg N$	$1.00 \lg N$	$1.00 \lg N$	$1.00 \lg N$	yes	<code>compareTo()</code>

Q. Can we do better?

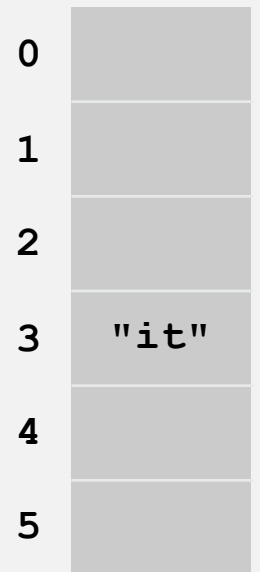
A. Yes, but with different access to the data.

# Hashing: basic plan

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

**Hash function.** Method for computing array index from key.

`hash("it") = 3`



0	
1	
2	
3	"it"
4	
5	

## Issues.

- Computing the hash function.
- Equality test: Method for checking whether two keys are equal.

# Hashing: basic plan

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

**Hash function.** Method for computing array index from key.

`hash("it") = 3`

0	
1	
2	
3	"it"
4	
5	

`hash("times") = 3` ??

## Issues.

- Computing the hash function.
- Equality test: Method for checking whether two keys are equal.
- Collision resolution: Algorithm and data structure to handle two keys that hash to the same array index.

## Classic space-time tradeoff.

- No space limitation: trivial hash function with key as index.
- No time limitation: trivial collision resolution with sequential search.
- Space and time limitations: hashing (the real world).

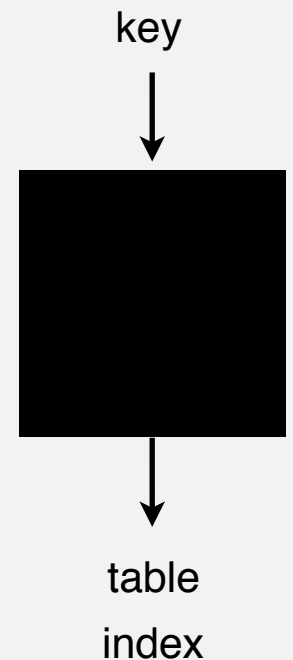
- ▶ **hash functions**
- ▶ separate chaining
- ▶ linear probing
- ▶ applications

# Computing the hash function

**Idealistic goal.** Scramble the keys uniformly to produce a table index.

- Efficiently computable.
- Each table index equally likely for each key.

thoroughly researched problem,  
still problematic in practical applications



**Ex 1. Phone numbers.**

- Bad: first three digits.
- Better: last three digits.

**Ex 2. Social Security numbers.**

- Bad: first three digits.
- Better: last three digits.

573 = California, 574 = Alaska  
(assigned in chronological order within geographic region)

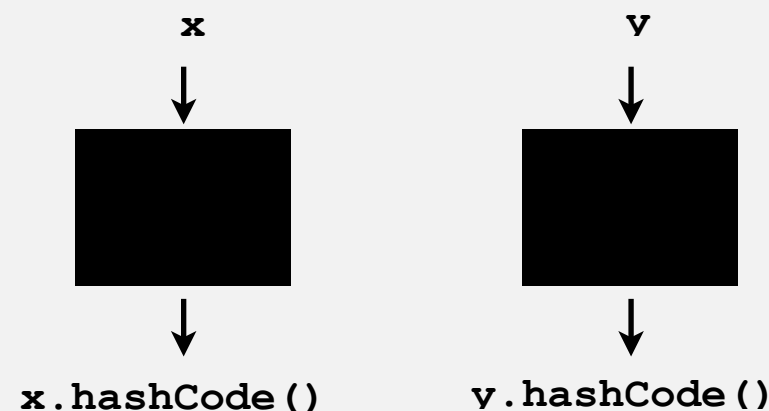
**Practical challenge.** Need different approach for each key type.

# Java's hash code 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())`.



**Default implementation.** Memory address of `x`.

**Trivial (but poor) implementation.** Always return 17.

**Customized implementations.** `Integer`, `Double`, `String`, `File`, `URL`, `Date`, ...

**User-defined types.** Users are on their own.



# Implementing hash code: integers, booleans, and doubles

```
public final class Integer
{
    private final int value;
    ...


    public int hashCode()
    { return value; }
}
```

```
public final class Boolean
{
    private final boolean value;
    ...

    public int hashCode()
    {
        if (value) return 1231;
        else      return 1237;
    }
}
```

```
public final class Double
{
    private final double value;
    ...

    public int hashCode()
    {
        long bits = doubleToLongBits(value);
        return (int) (bits ^ (bits >>> 32));
    }
}
```

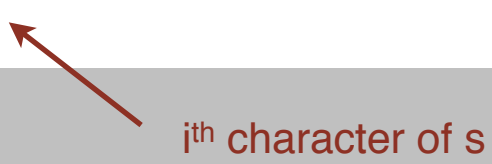


convert to IEEE 64-bit representation;  
xor most significant 32-bits  
with least significant 32-bits

# Implementing hash code: strings

```
public final class String
{
    private final char[] s;
    ...

    public int hashCode()
    {
        int hash = 0;
        for (int i = 0; i < length(); i++)
            hash = s[i] + (31 * hash);
        return hash;
    }
}
```



char	Unicode
...	...
'a'	97
'b'	98
'c'	99
...	...

- Horner's method to hash string of length  $L$ :  $L$  multiplies/adds.
- Equivalent to  $h = 31^{L-1} \cdot s^0 + \dots + 31^2 \cdot s^{L-3} + 31^1 \cdot s^{L-2} + 31^0 \cdot s^{L-1}$ .

Ex.

```
String s = "call";
int code = s.hashCode();
```


$$\begin{aligned} 3045982 &= 99 \cdot 31^3 + 97 \cdot 31^2 + 108 \cdot 31^1 + 108 \cdot 31^0 \\ &= 108 + 31 \cdot (108 + 31 \cdot (97 + 31 \cdot (99))) \end{aligned}$$

# War story: String hashing in Java

## String hashCode() in Java 1.1.

- For long strings: only examine 8-9 evenly spaced characters.
- Benefit: saves time in performing arithmetic.

```
public int hashCode()  
{  
    int hash = 0;  
    int skip = Math.max(1, length() / 8);  
    for (int i = 0; i < length(); i += skip)  
        hash = s[i] + (37 * hash);  
    return hash;  
}
```

- Downside: great potential for bad collision patterns.

# War story: String hashing in Java

## String hashCode() in Java 1.1.

- For long strings: only examine 8-9 evenly spaced characters.
- Benefit: saves time in performing arithmetic.

```
public int hashCode()  
{  
    int hash = 0;  
    int skip = Math.max(1, length() / 8);  
    for (int i = 0; i < length(); i += skip)  
        hash = s[i] + (37 * hash);  
    return hash;  
}
```

- Downside: great potential for bad collision patterns.

```
http://www.cs.princeton.edu/introcs/13loop/Hello.java  
http://www.cs.princeton.edu/introcs/13loop/Hello.class  
http://www.cs.princeton.edu/introcs/13loop/Hello.html  
http://www.cs.princeton.edu/introcs/12type/index.html  
↑      ↑      ↑      ↑      ↑      ↑      ↑      ↑
```

# Implementing hash code: user-defined types

```
public final class Transaction
{
    private final long who;
    private final Date when;
    private final String where;

    public Transaction(long who, Date when, String where)
    { /* as before */ }

    ...

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

```
public int hashCode()
{
    int hash = 17;
    hash = 31*hash + ((Long) who).hashCode();
    hash = 31*hash + when.hashCode();
    hash = 31*hash + where.hashCode();
    return hash;
}
```

nonzero constant

for primitive types, use  
`hashCode()`

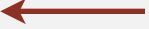
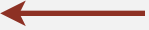
of wrapper type

for reference types,  
use `hashCode()`

typically a small prime

# Hash code design

"Standard" recipe for user-defined types.

- Combine each significant field using the  $31x + y$  rule.
- If field is a primitive type, use wrapper type `hashCode()`.
- If field is an array, apply to each element.  or use `Arrays.deepHashCode()`
- If field is a reference type, use `hashCode()`.  applies rule recursively

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

**In theory.** Need a theorem for each type to ensure reliability.

**Basic rule.** Need to use the whole key to compute hash code;  
consult an expert for state-of-the-art hash codes.

# 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

# 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



# 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  $-2^{31}$

# 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  $-2^{31}$

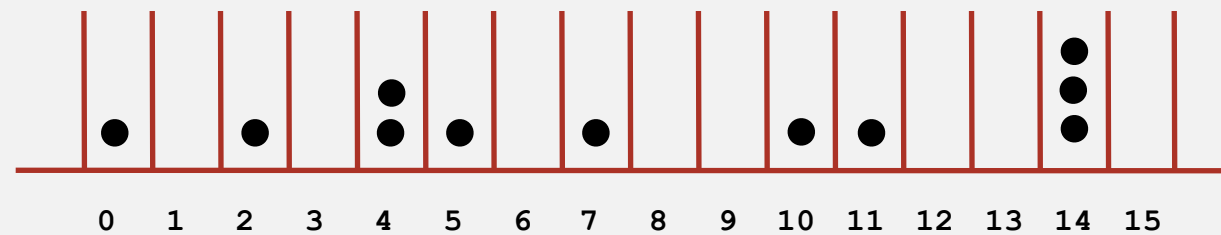
```
private int hash(Key key)
{   return (key.hashCode() & 0x7fffffff) % M;   }
```

correct

# Uniform hashing assumption

**Uniform hashing assumption.** Each key is equally likely to hash to an integer between 0 and  $M - 1$ .

**Bins and balls.** Throw balls uniformly at random into  $M$  bins.



**Birthday problem.** Expect two balls in the same bin after  $\sim \sqrt{\pi M / 2}$  tosses.

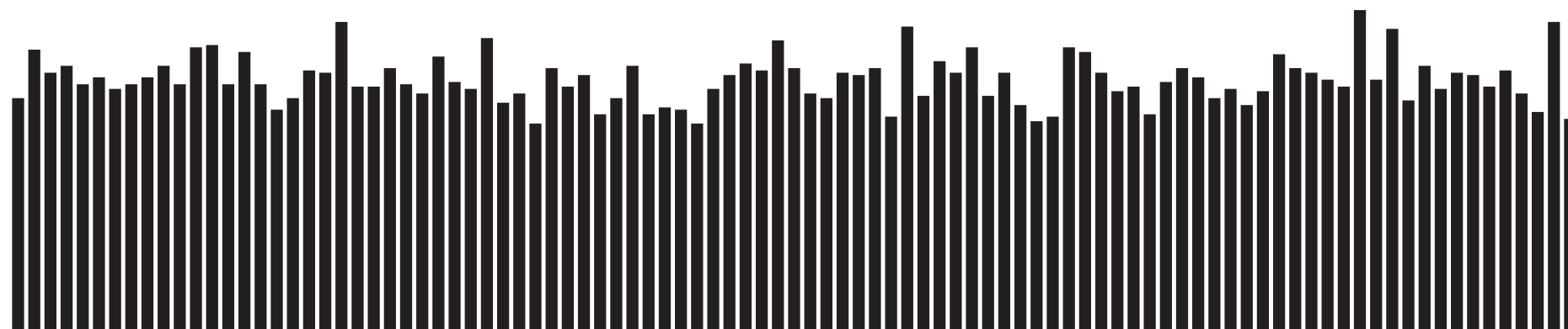
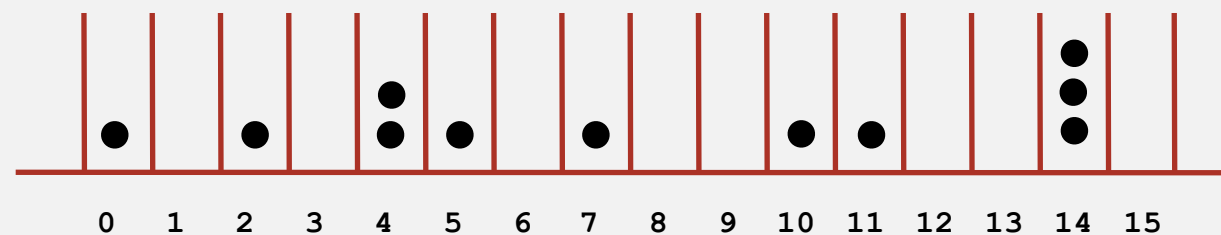
**Coupon collector.** Expect every bin has  $\geq 1$  ball after  $\sim M \ln M$  tosses.

**Load balancing.** After  $M$  tosses, expect most loaded bin has  $\Theta(\log M / \log \log M)$  balls.

# Uniform hashing assumption

**Uniform hashing assumption.** Each key is equally likely to hash to an integer between 0 and  $M - 1$ .

**Bins and balls.** Throw balls uniformly at random into  $M$  bins.



Hash value frequencies for words in Tale of Two Cities ( $M = 97$ )

Java's `String` data uniformly distribute the keys of Tale of Two Cities

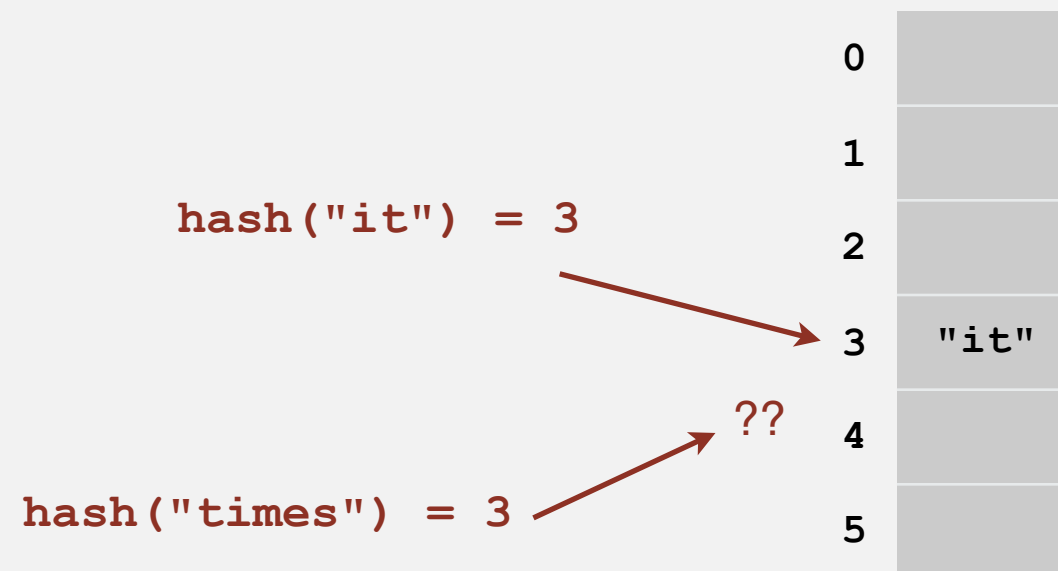
- ▶ hash functions
- ▶ **separate chaining**
- ▶ linear probing
- ▶ applications

# Collisions

**Collision.** Two distinct keys hashing to same index.

- Birthday problem  $\Rightarrow$  can't avoid collisions unless you have a ridiculous (quadratic) amount of memory.
- Coupon collector + load balancing  $\Rightarrow$  collisions will be evenly distributed.

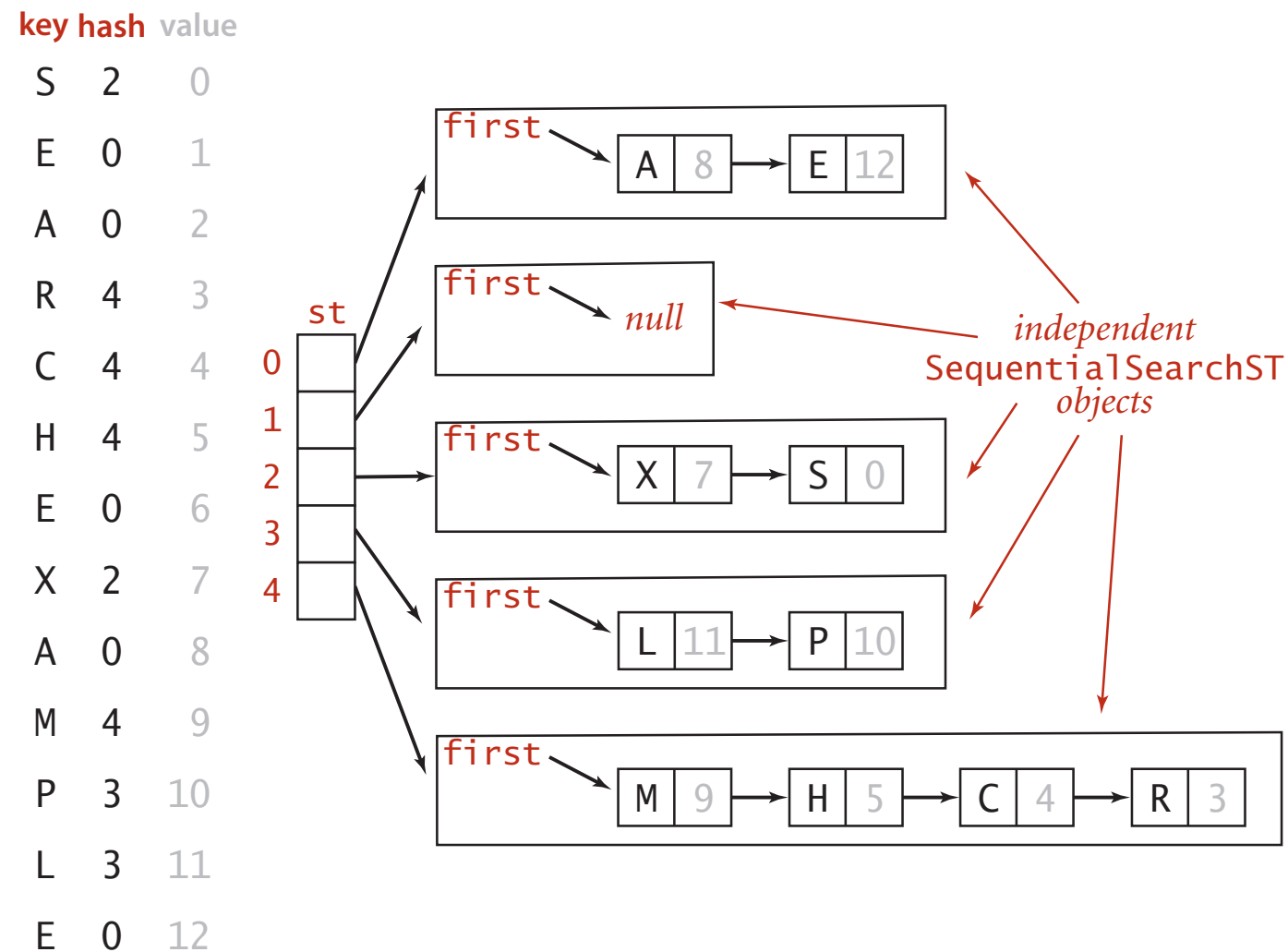
**Challenge.** Deal with collisions efficiently.



# Separate chaining ST

Use an array of  $M < N$  linked lists. [H. P. Luhn, IBM 1953]

- Hash: map key to integer  $i$  between 0 and  $M - 1$ .
- Insert: put at front of  $i^{\text{th}}$  chain (if not already there).
- Search: only need to search  $i^{\text{th}}$  chain.



Hashing with separate chaining for standard indexing client

# Separate chaining ST: Java implementation

```
public class SeparateChainingHashST<Key, Value>
{
    private int N;        // number of key-value pairs
    private int M;        // hash table size
    private SequentialSearchST<Key, Value> [] st;    // array of STs

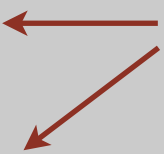
    public SeparateChainingHashST()
    { this(997); }

    public SeparateChainingHashST(int M)
    {
        this.M = M;
        st = (SequentialSearchST<Key, Value>[]) new SequentialSearchST[M];
        for (int i = 0; i < M; i++)
            st[i] = new SequentialSearchST<Key, Value>();
    }

    private int hash(Key key)
    { return (key.hashCode() & 0x7fffffff) % M; }

    public Value get(Key key)
    { return st[hash(key)].get(key); }

    public void put(Key key, Value val)
    { st[hash(key)].put(key, val); }
}
```



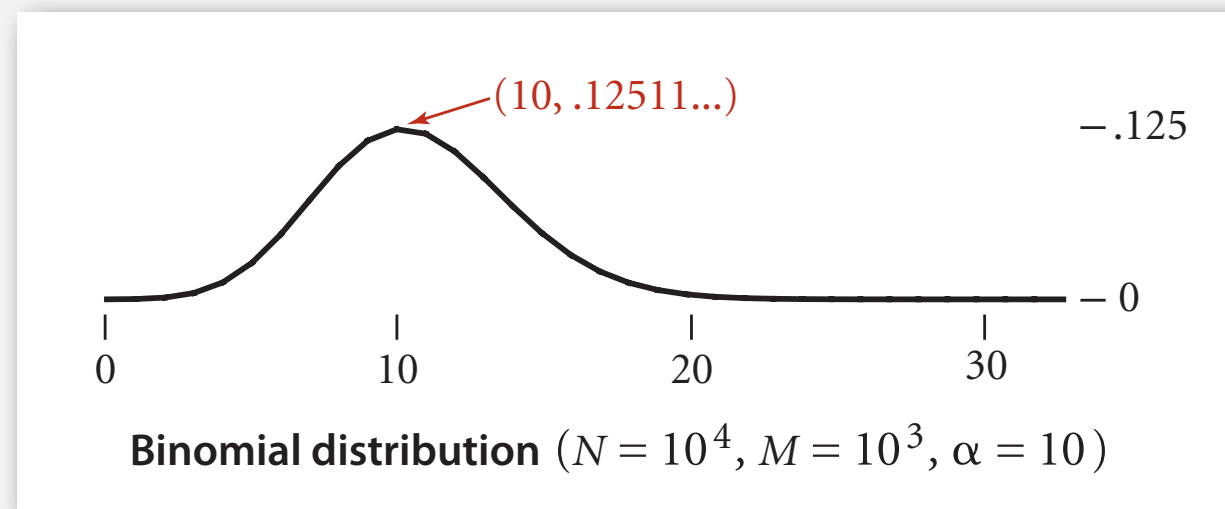
array doubling and halving code omitted



# Analysis of separate chaining

**Proposition.** Under uniform hashing assumption, probability that the number of keys in a list is within a constant factor of  $N/M$  is extremely close to 1.

**Pf sketch.** Distribution of list size obeys a binomial distribution.



**Consequence.** Number of probes for search/insert is proportional to  $N/M$ .

- $M$  too large  $\Rightarrow$  too many empty chains.
- $M$  too small  $\Rightarrow$  chains too long.
- Typical choice:  $M \sim N/5 \Rightarrow$  constant-time ops.

↑  
M times faster than  
sequential search

# ST implementations: summary

implementation	guarantee			average case			ordered iteration?	operations on keys
	search	insert	delete	search hit	insert	delete		
sequential search (linked list)	N	N	N	N/2	N	N/2	no	<b>equals()</b>
binary search (ordered array)	lg N	N	N	lg N	N/2	N/2	yes	<b>compareTo()</b>
BST	N	N	N	1.38 lg N	1.38 lg N	?	yes	<b>compareTo()</b>
red-black tree	2 lg N	2 lg N	2 lg N	1.00 lg N	1.00 lg N	1.00 lg N	yes	<b>compareTo()</b>
separate chaining	lg N *	lg N *	lg N *	3-5 *	3-5 *	3-5 *	no	<b>equals()</b>

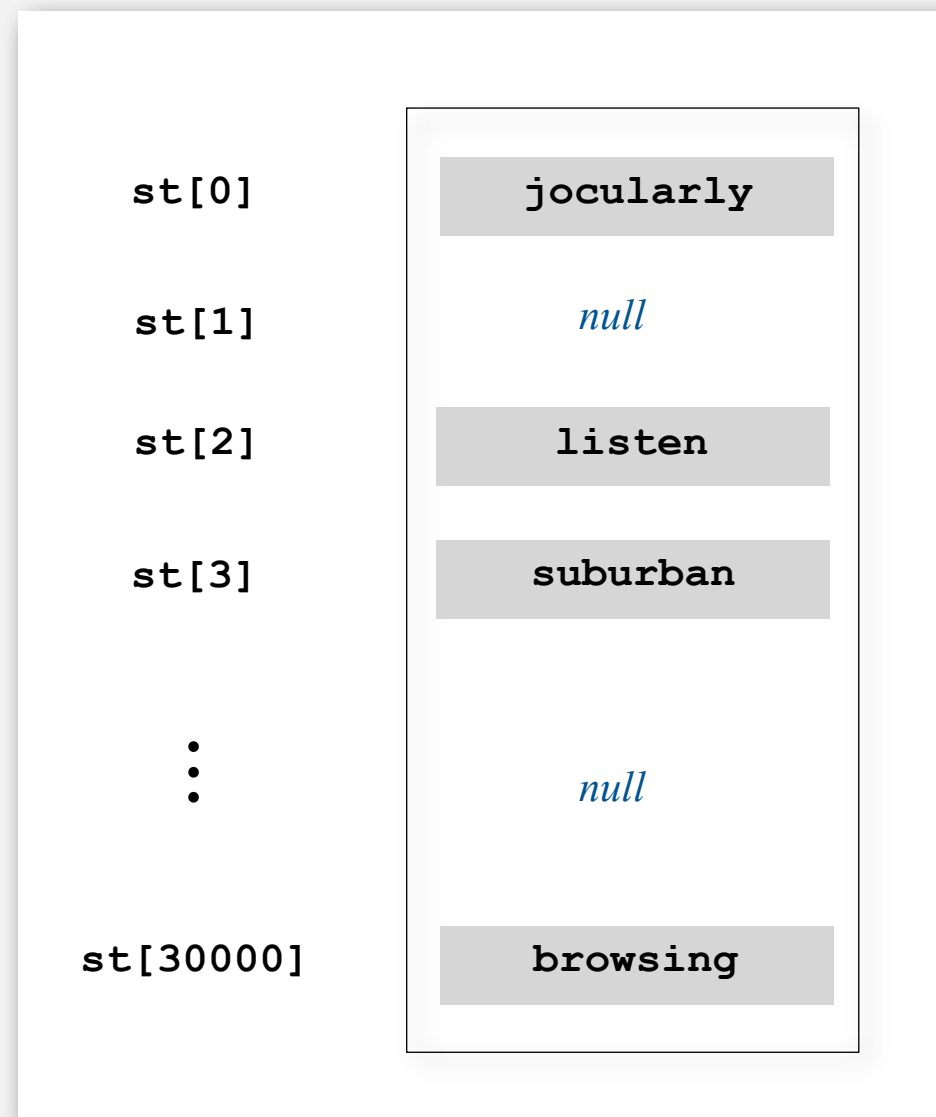
\* under uniform hashing assumption

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

# Collision resolution: open addressing

Open addressing. [Amdahl-Boehme-Rocherster-Samuel, IBM 1953]

When a new key collides, find next empty slot, and put it there.



linear probing ( $M = 30001$ ,  $N = 15000$ )

# Linear probing

Use an array of size  $M > N$ .

- Hash: map key to integer  $i$  between 0 and  $M - 1$ .
- Insert: put at table index  $i$  if free; if not try  $i + 1$ ,  $i + 2$ , etc.
- Search: search table index  $i$ ; if occupied but no match, try  $i + 1$ ,  $i + 2$ , etc.

-	-	-	S	H	-	-	A	C	E	R	-	-
0	1	2	3	4	5	6	7	8	9	10	11	12

# Linear probing

Use an array of size  $M > N$ .

- Hash: map key to integer  $i$  between 0 and  $M - 1$ .
- Insert: put at table index  $i$  if free; if not try  $i + 1$ ,  $i + 2$ , etc.
- Search: search table index  $i$ ; if occupied but no match, try  $i + 1$ ,  $i + 2$ , etc.

-	-	-	S	H	-	-	A	C	E	R	-	-
0	1	2	3	4	5	6	7	8	9	10	11	12

-	-	-	S	H	-	-	A	C	E	R	I	-
0	1	2	3	4	5	6	7	8	9	10	11	12

insert I  
hash(I) = 11

# Linear probing

Use an array of size  $M > N$ .

- Hash: map key to integer  $i$  between 0 and  $M - 1$ .
- Insert: put at table index  $i$  if free; if not try  $i + 1$ ,  $i + 2$ , etc.
- Search: search table index  $i$ ; if occupied but no match, try  $i + 1$ ,  $i + 2$ , etc.

-	-	-	S	H	-	-	A	C	E	R	-	-
0	1	2	3	4	5	6	7	8	9	10	11	12

-	-	-	S	H	-	-	A	C	E	R	I	-
0	1	2	3	4	5	6	7	8	9	10	11	12

insert I  
hash(I) = 11

-	-	-	S	H	-	-	A	C	E	R	I	N
0	1	2	3	4	5	6	7	8	9	10	11	12

insert N  
hash(N) = 8

# Linear probing: trace of standard indexing client

key	hash	value	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
S	6	0							S									
									0									
E	10	1							S				E					
									0				1					
A	4	2					A		S				E					
							2		0				1					
R	14	3					A		S				E				R	
							2		0				1				3	
C	5	4					A	C	S				E				R	
							2	5	0				1				3	
H	4	5					A	C	S	H			E				R	
							2	5	0	5			1				3	
E	10	6					A	C	S	H			E				R	
							2	5	0	5			6				3	
X	15	7					A	C	S	H			E				R	X
							2	5	0	5			6				3	7
A	4	8					A	C	S	H			E				R	X
							8	5	0	5			6				3	7
M	1	9		M			A	C	S	H			E				R	X
				9			8	5	0	5			6				3	7
P	14	10	P	M			A	C	S	H			E				R	X
			10	9			8	5	0	5			6				3	7
L	6	11	P	M			A	C	S	H	L		E				R	X
			10	9			8	5	0	5	11		6				3	7
E	10	12	P	M			A	C	S	H	L		E				R	X
			10	9			8	5	0	5	11		12				3	7

entries in red are new

entries in gray are untouched

keys in black are probes

probe sequence wraps to 0

keys[]

vals[]



# Linear probing ST implementation


```
public class LinearProbingHashST<Key, Value>
{
    private int M = 30001;
    private Value[] vals = (Value[]) new Object[M];
    private Key[] keys = (Key[]) new Object[M];

    private int hash(Key key) { /* as before */ }

    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;
    }

    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;
    }
}
```

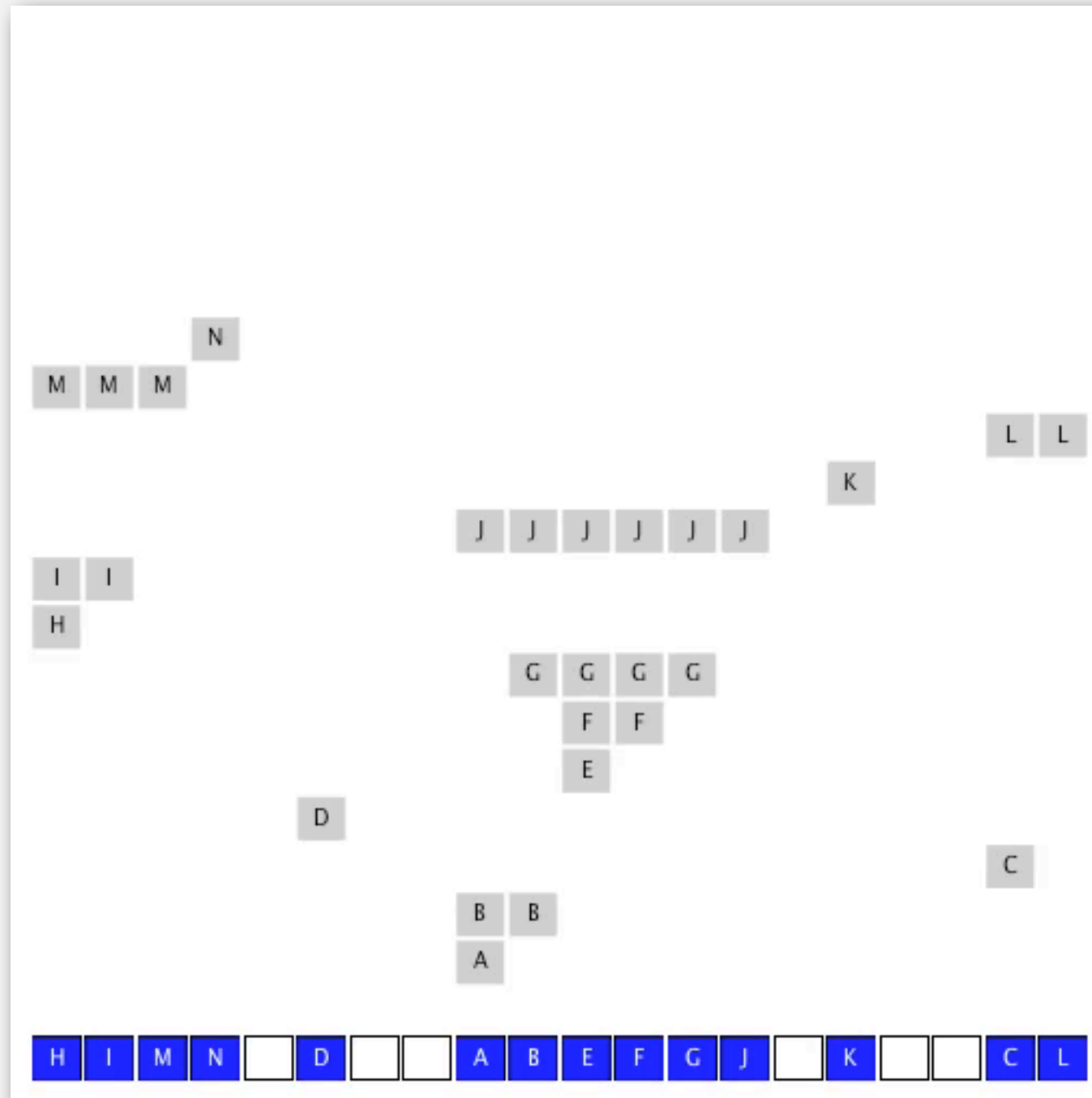
array doubling  
and halving  
code omitted



# Clustering

**Cluster.** A contiguous block of items.

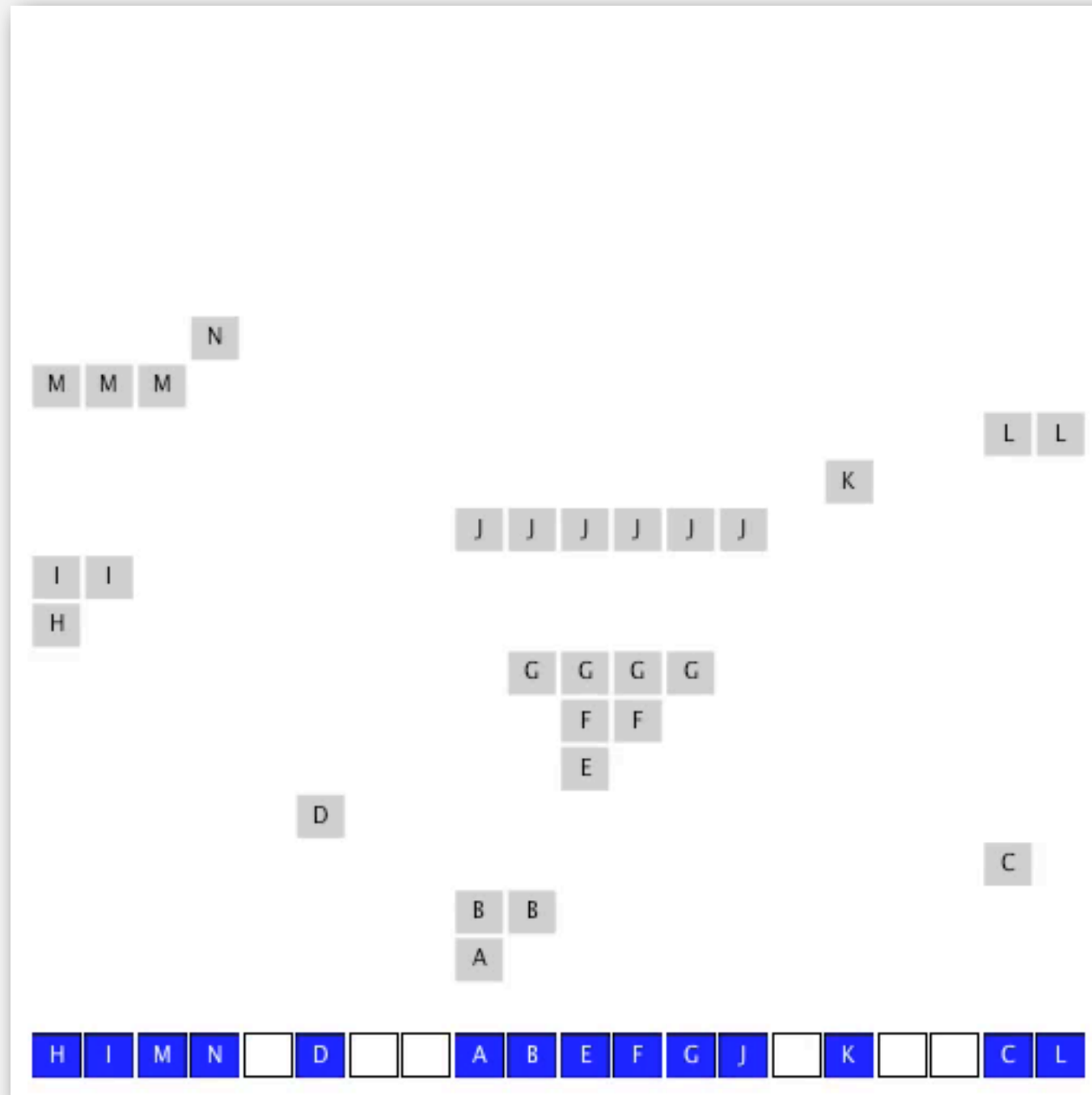
**Observation.** New keys likely to hash into middle of big clusters.



# Clustering

**Cluster.** A contiguous block of items.

**Observation.** New keys likely to hash into middle of big clusters.

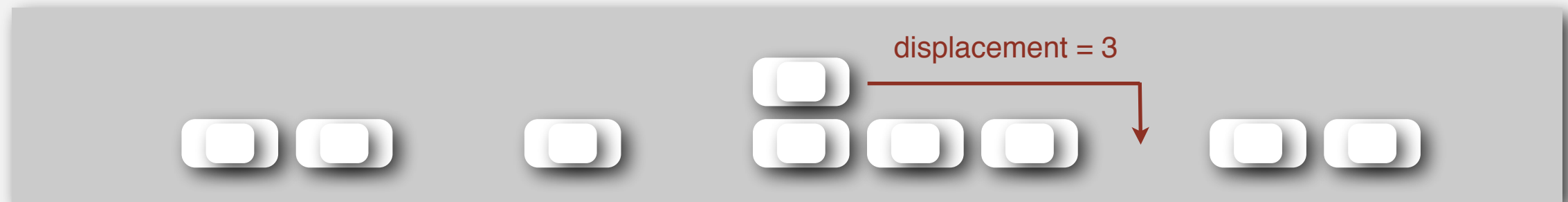


# Knuth's parking problem

**Model.** Cars arrive at one-way street with  $M$  parking spaces.

Each desires a random space  $i$ : if space  $i$  is taken, try  $i + 1, i + 2$ , etc.

**Q.** What is mean displacement of a car?



**Half-full.** With  $M / 2$  cars, mean displacement is  $\sim 3 / 2$ .

**Full.** With  $M$  cars, mean displacement is  $\sim \sqrt{\pi M / 8}$

# Analysis of linear probing

**Proposition.** Under uniform hashing assumption, the average number of probes in a hash table of size  $M$  that contains  $N = \alpha M$  keys is:

$$\begin{array}{cc} \sim \frac{1}{2} \left( 1 + \frac{1}{1 - \alpha} \right) & \sim \frac{1}{2} \left( 1 + \frac{1}{(1 - \alpha)^2} \right) \\ \text{search hit} & \text{search miss / insert} \end{array}$$

**Pf.** [Knuth 1962] A landmark in analysis of algorithms.

## 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  $5/2$

# ST implementations: summary

implementation	guarantee			average case			ordered iteration?	operations on keys
	search	insert	delete	search hit	insert	delete		
sequential search (linked list)	N	N	N	N/2	N	N/2	no	<b>equals()</b>
binary search (ordered array)	lg N	N	N	lg N	N/2	N/2	yes	<b>compareTo()</b>
BST	N	N	N	1.38 lg N	1.38 lg N	?	yes	<b>compareTo()</b>
red-black tree	2 lg N	2 lg N	2 lg N	1.00 lg N	1.00 lg N	1.00 lg N	yes	<b>compareTo()</b>
separate chaining	lg N *	lg N *	lg N *	3-5 *	3-5 *	3-5 *	no	<b>equals()</b>
linear probing	lg N *	lg N *	lg N *	3-5 *	3-5 *	3-5 *	no	<b>equals()</b>

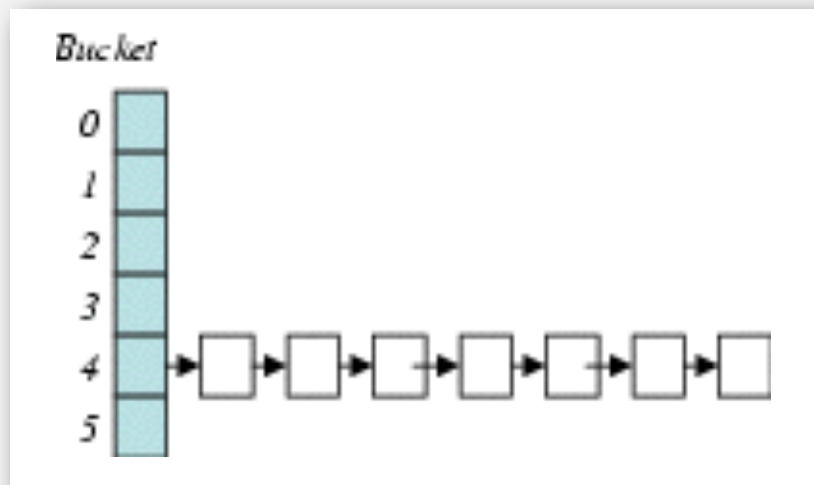
\* under uniform hashing assumption

# War story: algorithmic complexity attacks

Q. Is the uniform hashing assumption important in practice?

A. Obvious situations: aircraft control, nuclear reactor, pacemaker.

A. Surprising situations: **denial-of-service** attacks.



malicious adversary learns your hash function  
(e.g., by reading Java API) and causes a big pile-up  
in single slot that grinds performance to a halt

**Real-world exploits.** [Crosby-Wallach 2003]

- Bro server: send carefully chosen packets to DOS the server, using less bandwidth than a dial-up modem.
- Perl 5.8.0: insert carefully chosen strings into associative array.
- Linux 2.4.20 kernel: save files with carefully chosen names.

# Algorithmic complexity attack on Java

**Goal.** Find family of strings with the same hash code.

**Solution.** The base-31 hash code is part of Java's string API.

key	hashCode ()
"Aa"	2112
"BB"	2112



# Algorithmic complexity attack on Java

**Goal.** Find family of strings with the same hash code.

**Solution.** The base-31 hash code is part of Java's string API.

key	hashCode ()
"Aa"	2112
"BB"	2112

key	hashCode ()
"AaAaAaAa"	-540425984
"AaAaAaBB"	-540425984
"AaAaBBaA"	-540425984
"AaAaBBBB"	-540425984
"AaBBaAaA"	-540425984
"AaBBaABB"	-540425984
"AaBBBBaA"	-540425984
"AaBBBBBB"	-540425984

key	hashCode ()
"BBaAaAaA"	-540425984
"BBaAaBB"	-540425984
"BBaBBaA"	-540425984
"BBaBBBB"	-540425984
"BBBBaAaA"	-540425984
"BBBBaABB"	-540425984
"BBBBBBaA"	-540425984
"BBBBBBBB"	-540425984

**$2^N$  strings of length  $2N$  that hash to same value!**

# Diversion: one-way hash functions

**One-way hash function.** "Hard" to find a key that will hash to a desired value (or two keys that hash to same value).

**Ex.** MD4, MD5, SHA-0, SHA-1, SHA-2, WHIRLPOOL, RIPEMD-160, ...

known to be insecure

```
String password = args[0];  
MessageDigest sha1 = MessageDigest.getInstance("SHA1");  
byte[] bytes = sha1.digest(password);  
  
/* prints bytes as hex string */
```

**Applications.** Digital fingerprint, message digest, storing passwords.

**Caveat.** Too expensive for use in ST implementations.

# Separate chaining vs. linear probing

## Separate chaining.

- Easier to implement delete.
- Performance degrades gracefully.
- Clustering less sensitive to poorly-designed hash function.

## Linear probing.

- Less wasted space.
- Better cache performance.

# Hashing: variations on the theme

Many improved versions have been studied.

**Two-probe hashing.** (separate-chaining variant)

- Hash to two positions, put key in shorter of the two chains.
- Reduces expected length of the longest chain to  $\log \log N$ .

**Double hashing.** (linear-probing variant)

- Use linear probing, but skip a variable amount, not just 1 each time.
- Effectively eliminates clustering.
- Can allow table to become nearly full.
- Difficult to implement delete.

# Hashing vs. balanced search trees

## Hashing.

- Simpler to code.
- No effective alternative for unordered keys.
- Faster for simple keys (a few arithmetic ops versus  $\log N$  compares).
- Better system support in Java for strings (e.g., cached hash code).

## Balanced search trees.

- Stronger performance guarantee.
- Support for ordered ST operations.
- Easier to implement `compareTo()` correctly than `equals()` and `hashCode()`.

## Java system includes both.

- Red-black trees: `java.util.TreeMap`, `java.util.TreeSet`.
- Hashing: `java.util.HashMap`, `java.util.IdentityHashMap`.

- ▶ sets
- ▶ dictionary clients
- ▶ indexing clients
- ▶ sparse vectors
- ▶ **challenges**

# ST implementations: summary

implementation	guarantee			average case			ordered iteration?	operations on keys
	search	insert	delete	search hit	insert	delete		
sequential search (linked list)	N	N	N	N/2	N	N/2	no	<b>equals()</b>
binary search (ordered array)	lg N	N	N	lg N	N/2	N/2	yes	<b>compareTo()</b>
BST	N	N	N	1.38 lg N	1.38 lg N	?	yes	<b>compareTo()</b>
red-black tree	2 lg N	2 lg N	2 lg N	1.00 lg N	1.00 lg N	1.00 lg N	yes	<b>compareTo()</b>
separate chaining	lg N *	lg N *	lg N *	3-5 *	3-5 *	3-5 *	no	<b>equals()</b>
linear probing	lg N *	lg N *	lg N *	3-5 *	3-5 *	3-5 *	no	<b>equals()</b>

\* under uniform hashing assumption

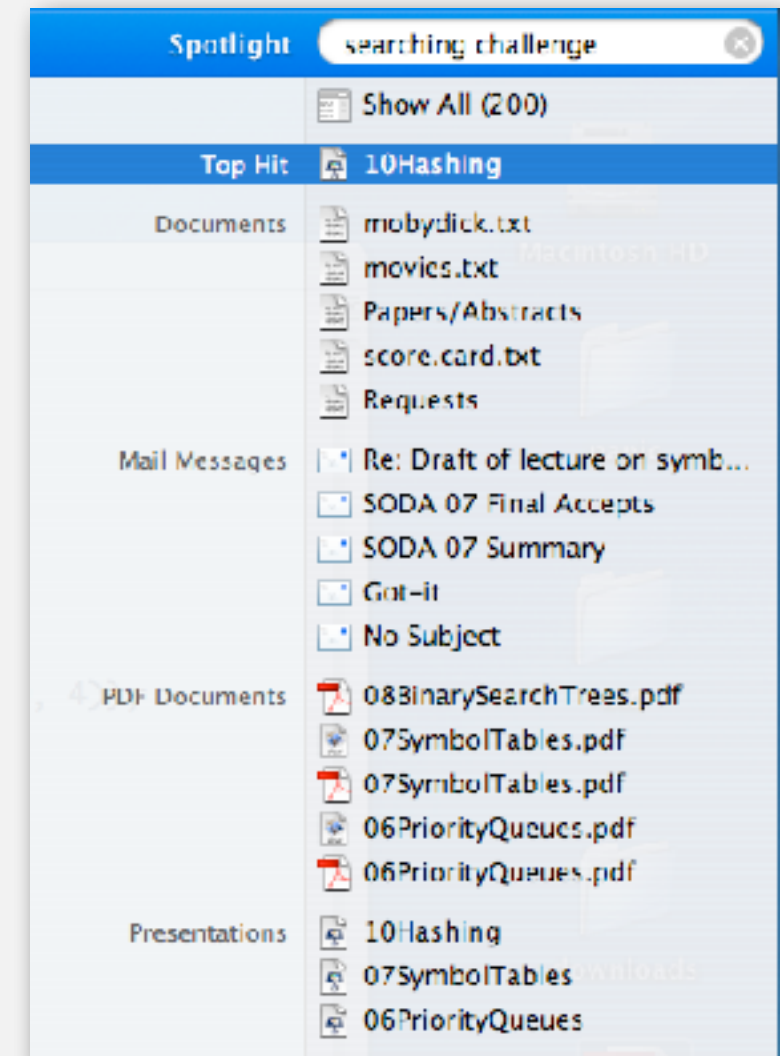
# Searching challenge 1

**Problem.** Index for a PC or the web.

**Assumptions.** 1 billion++ words to index.

**Which searching method to use?**

- Hashing
- Red-black-trees
- Doesn't matter much.





# Searching challenge 1

**Problem.** Index for a PC or the web.

**Assumptions.** 1 billion++ words to index.

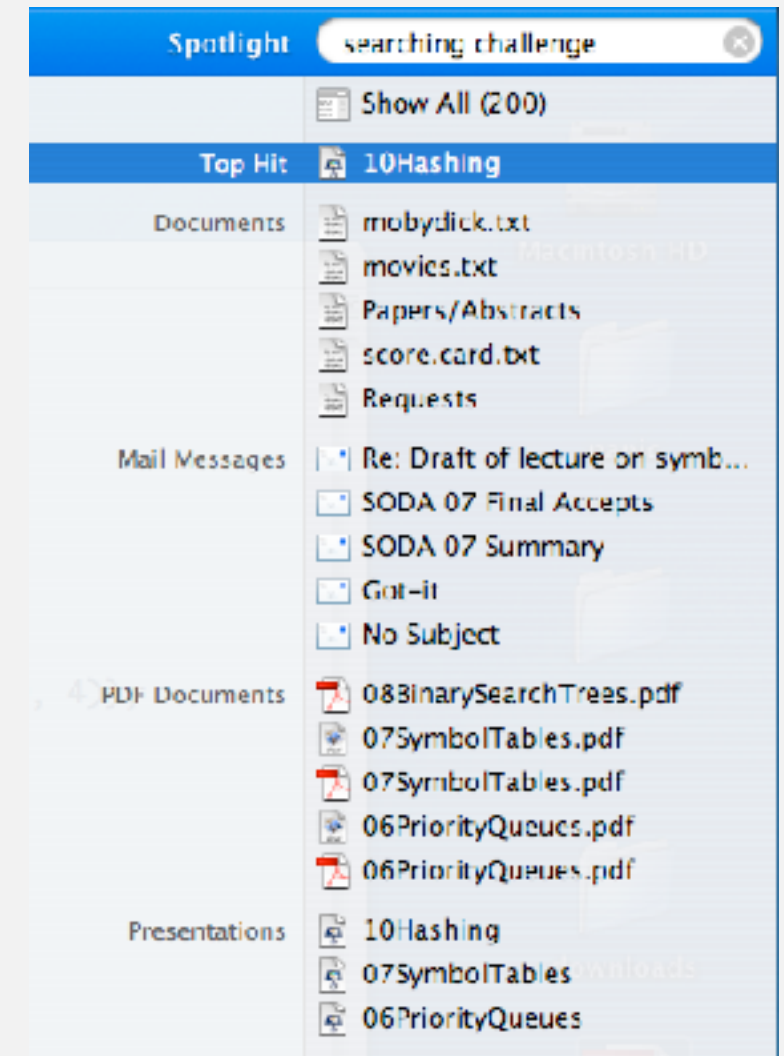
**Which searching method to use?**

- ✓ Hashing
- Red-black-trees ← too much space
- Doesn't matter much.

**Solution.** Symbol table with:

- Key = query string.
- Value = set of pointers to files.

← sort the (relatively few) search hits



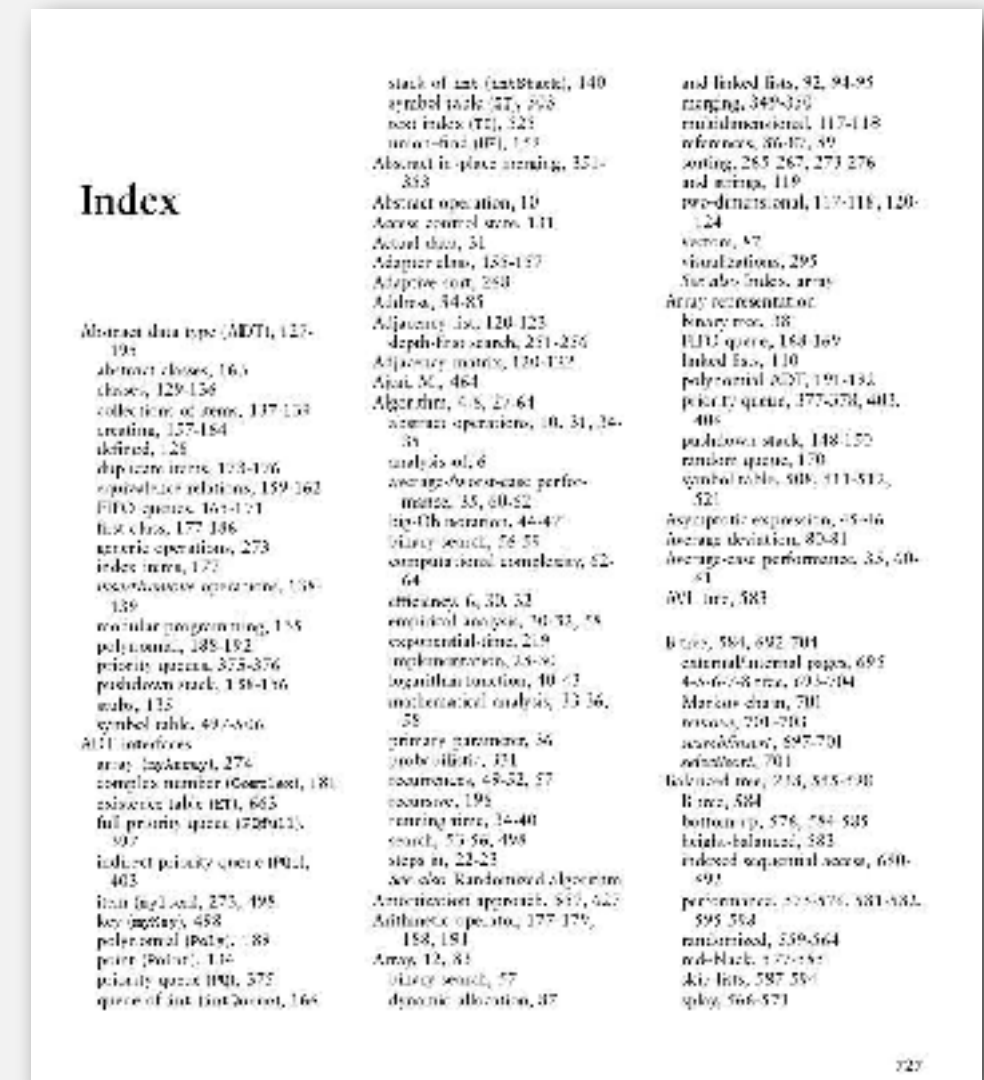
## Searching challenge 2

## Problem. Index for an e-book.

**Assumptions.** Book has 100,000+ words.

## Which searching method to use?

1. Hashing
2. Red-black-tree
3. Doesn't matter much.



## Searching challenge 2

## Problem. Index for an e-book.

**Assumptions.** Book has 100,000+ words.

## Which searching method to use?

1. Hashing
- ✓ 2. Red-black-tree ← need ordered iteration
3. Doesn't matter much.

**Solution.** Symbol table with:

- Key = index term.
- Value = ordered set of pages on which term appears.

# Index

Abstract data type (ADT), 127–139  
  abstract classes, 165  
  classes, 129–136  
  collection of items, 137–138  
  creating, 127–134  
  defined, 128  
  dynamic items, 133–136  
  equivalence relations, 159–162  
  FIFO queue, 163–171  
  list class, 177–186  
  numeric operations, 273  
  index items, 173  
  recurrence operations, 158–159  
  recursive programming, 135  
  polynomial, 188–192  
  priority queue, 373–376  
  pushdown stack, 158–159  
  sets, 135  
  symbol table, 497–506

ADT interfaces  
  array, *bigArray*, 274  
  complex number *Complex*, 181  
  existence table *HT*, 663  
  full priority queue *FullPQ*, 397  
  indexed priority queue *IPQ*, 403  
  list (*bigList*), 275, 295  
  key (*bigKey*), 458  
  polynomial (*Pol*), 183  
  point (*Point*), 134  
  priority queue *PQ*, 375  
  queue of int (*IntQueue*), 166  
  stack of int (*IntStack*), 140  
  symbol table (*ST*), 508  
  text index (*TI*), 523  
  union-find (*UF*), 153

Abstract in-place sorting, 351–353

Abstract operation, 19

Access control items, 111

Actual data, 31

Adapter class, 183–187

Adaptive sort, 268

Address, 34–35

Adjacency list, 120–123  
  depth-first search, 231–235

Adjacency matrix, 120–122

Alan, M., 464

Algorithm, 6–8, 27–64  
  numeric operations, 10, 31, 34–35  
  analysis of, 6  
  average-case-case performance, 35, 60–62  
  big-Oh notation, 44–47  
  binary search, 56–55  
  computational modeling, 62–64  
  efficiency, 6, 30, 32  
  empirical analysis, 30–32, 58  
  exponential-time, 219  
  implementation, 28–30  
  logarithmic function, 40–42  
  mathematical analysis, 33–36, 55  
  primitive grammar, 56  
  proofs of, 33  
  recurrences, 48–52, 57  
  recursion, 195  
  running time, 34–40  
  search, 53–56, 498  
  steps in, 22–23  
  see also Randomized Algorithm

Amortization approach, 510, 522

Authentic operators, 177–179, 184, 181

Array, 12, 81  
  binary search, 57  
  dynamic allocation, 37  
  and linked lists, 92, 94–95  
  merging, 349–350  
  multidimensional, 117–118  
  references, 36–37, 39  
  sorting, 283–287, 273–276  
  and strings, 119  
  two-dimensional, 117–118, 120–124  
  vector, 87  
  visualizations, 295  
  see also Index, array

Array representation  
  binary max, 38  
  FIFO queue, 164–169  
  linked list, 110  
  polynomial ADT, 191–192  
  priority queue, 377–378, 401, 408  
  pushdown stack, 148–150  
  random queue, 170  
  symbol table, 508, 511–513, 521

Asymptotic expression, 45–46

Average deviation, 82–81

Average-case performance, 35, 60–61

AVL tree, 583

Baker, 584, 592–704

External internal pages, 693  
  4–6–6–6–6 tree, 693–694  
  Marion diagram, 701  
  reversal, 704–704  
  searching, 597–701  
  selection, 701

Interval map, 214, 515–516

Items, 584  
  bottom up, 576, 584–585  
  heap-balanced, 583  
  indexed sequential access, 690–691

Performance, 575–576, 581–582, 595–596  
  normalized, 559–564  
  red-black, 577–581  
  skip lists, 587–591  
  splay, 566–571

727