http://algs4.cs.princeton.edu

## 3.4 HASH TABLES

▸ *hash functions*

▸ *separate chaining*

▸ *linear probing*

▸ *context*

# Symbol table implementations:  summary

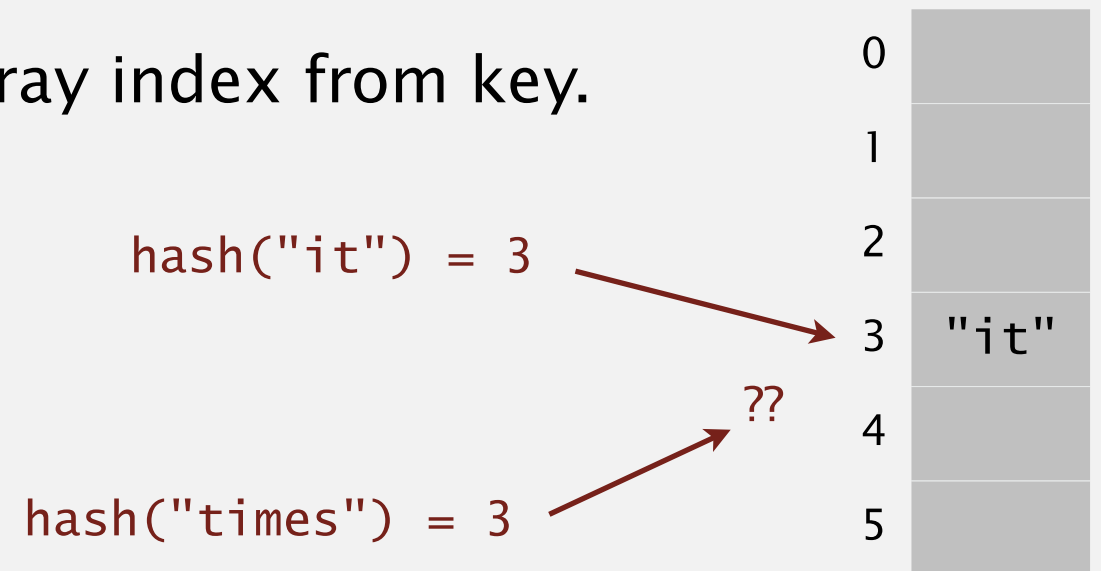| implementation | guarantee | | | average case | | | ordered ops? | key interface |
|---|---|---|---|---|---|---|---|---|
| | search | insert | delete | search hit | insert | delete | | |
| **sequential search (unordered list)** | $N$ | $N$ | $N$ | ½ $N$ | $N$ | ½ $N$ | | `equals()` |
| **binary search (ordered array)** | lg $N$ | $N$ | $N$ | lg $N$ | ½ $N$ | ½ $N$ | ✔ | `compareTo()` |
| **BST** | $N$ | $N$ | $N$ | 1.39 lg $N$ | 1.39 lg $N$ | $\sqrt{N}$ | ✔ | `compareTo()` |
| **red–black BST** | 2 lg $N$ | 2 lg $N$ | 2 lg $N$ | 1.0 lg $N$ | 1.0 lg $N$ | 1.0 lg $N$ | ✔ | `compareTo()` |

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
```

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

| | |
|---|---|
| 0 | |
| 1 | |
| 2 | |
| 3 | "it" |
| 4 | ?? |
| 5 | |

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).

# 3.4 HASH TABLES

▸ *hash functions*

▸ *separate chaining*

▸ *linear probing*

▸ *context*

Algorithms

ROBERT SEDGEWICK | KEVIN WAYNE

http://algs4.cs.princeton.edu

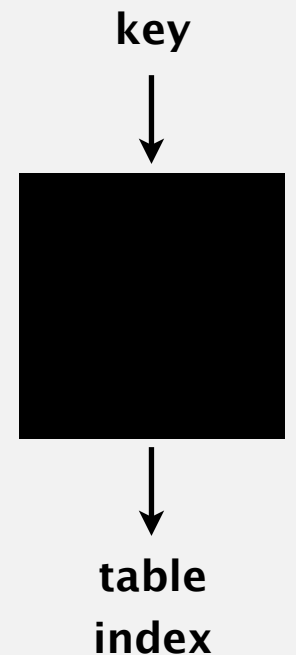# Computing the hash function

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*

**key**

**table
index**

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
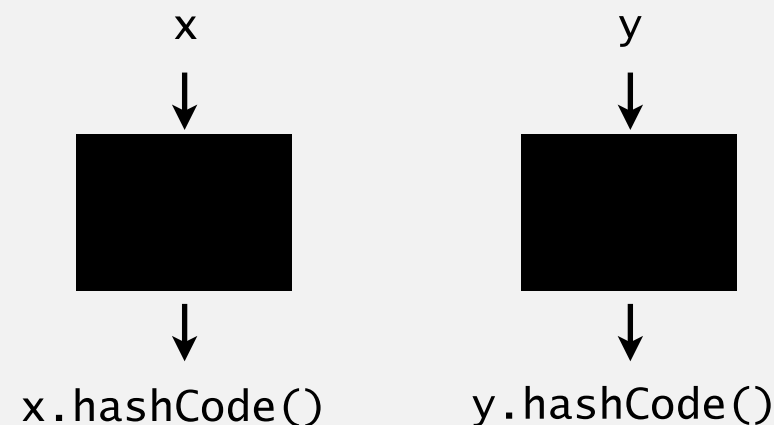(assigned in chronological order within geographic region)

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())`.



```
x                    y
↓                    ↓
[■■■■]            [■■■■]
↓                    ↓
x.hashCode()     y.hashCode()
```

Default implementation.  Memory address of `x`.

Legal (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

**Java library implementations**

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

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

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

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

```java
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

Warning: -0.0 and +0.0 have different hash codes

# Implementing hash code:  strings

**Java library implementation**

```
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;
   }
}
```

ith character of s

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

- Horner's method to hash string of length $L$:  $L$ multiplies/adds.

- Equivalent to  $h = s[0] \cdot 31^{L-1} + \ldots + s[L-3] \cdot 31^2 + s[L-2] \cdot 31^1 + s[L-1] \cdot 31^0$.

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

$3045982 = 99 \cdot 31^3 + 97 \cdot 31^2 + 108 \cdot 31^1 + 108 \cdot 31^0$
$= 108 + 31 \cdot (108 + 31 \cdot (97 + 31 \cdot (99)))$
(Horner's method)

# Implementing hash code:  strings

Performance optimization.

- Cache the hash value in an instance variable.
- Return cached value.

```
public final class String
{
    private int hash = 0;              ←——— cache of hash code
    private final char[] s;

    ...

    public int hashCode()
    {
        int h = hash;                  ←——— return cached value
        if (h != 0) return h;
        for (int i = 0; i < length(); i++)
            h = s[i] + (31 * h);
        hash = h;                      ←——— store cache of hash code
        return h;
    }
}
```

Q.  What if `hashCode()` of string is 0?

# Implementing hash code:  user-defined types

```
public final class Transaction implements Comparable<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 = 17;
      hash = 31*hash + who.hashCode();
      hash = 31*hash + when.hashCode();
      hash = 31*hash + ((Double) amount).hashCode();
      return hash;
   }
}
```

nonzero constant

for reference types,
use hashCode()

for primitive types,
use hashCode()
of wrapper type

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 `null`, return $0$.
- If field is a reference type, use `hashCode()`. ⟵ applies rule recursively
- If field is an array, apply to each entry. ⟵ or use `Arrays.deepHashCode()`

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

In theory.  Keys are bitstring; "universal" hash functions exist.

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

x

↓

↓

x.hashCode()

↓

```
private int hash(Key key)
{   return key.hashCode() % M;   }
```
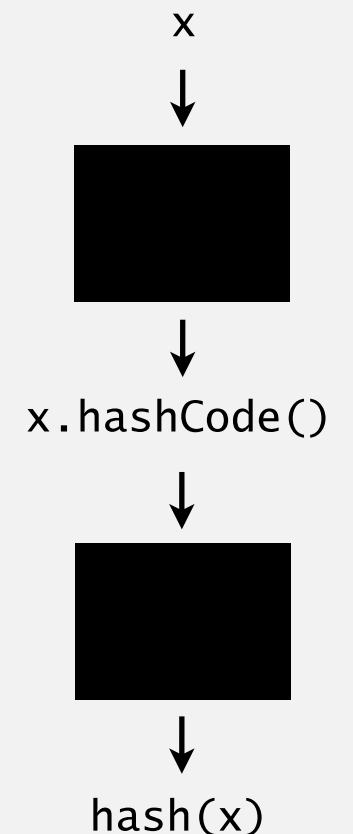
**bug**

hash(x)

```
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;   }
```
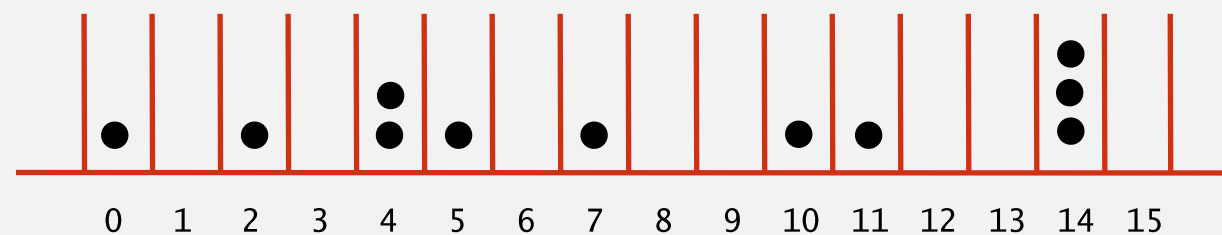
**correct**

# 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.



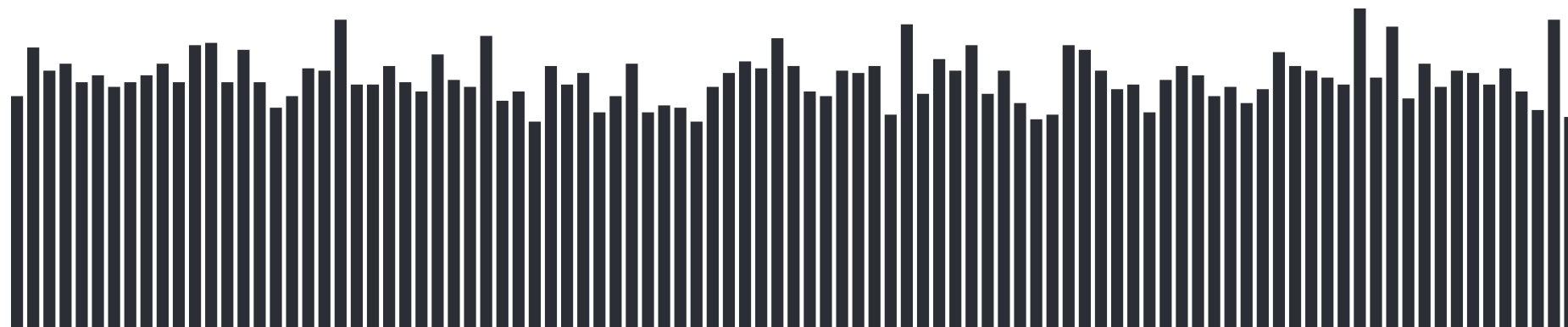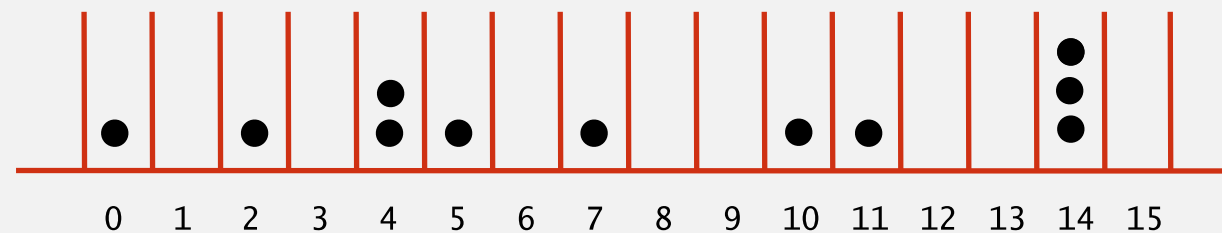Birthday problem. Expect two balls in the same bin after $\sim \sqrt{\pi M / 2}$ tosses.

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