# 6.006 Recitation

Build 2008.10

# Coming Up Next...

- Hashing in theory and in Python

- Bad hash functions

- Mutable dictionary keys

- Hashes for basic data types in Python

# Why Hashing

- Useless from a theoretical standpoint
  - $O(N)$ / op worst-case, not fit for proofs
- Used everywhere (dictionaries, indices)
  - $O(1)$ / op is smokin' hot / fast
  - Simple - small constant factor
  - Relies on black magic

# Hashing pwns BSTs?

- BSTs

  - O(lg(N)) / op

  - guaranteed upper bound (worst-case)

  - comparison model (an order relation on keys is sufficient)

  - pwns in real-time

- Hashing

  - O(1) / op avg-case

  - no guarantees for worst-case -- O(N)

  - intimate knowledge of keys (via magic inside the hash function)

  - rocks for most cases

# Real Life Hashing I

- Application: Keeping library cards
  - 4x6" card for each book
  - filing by the 1st letter of the book title
    - e.g. "Differential Equations" goes to D
  - no sorting asides from mechanism above

# Real Life Hashing II

- **filing** is uncool, let's think of **bucketing**
  - 26 buckets, labeled 'A' - 'Z'
  - Books are bucketed by 1st letter in title
  - Time to find a book ~ bucket size

# Real Life Hashing III

- What sucks in the scheme above?
  - Common prefixes
    - "The ...", "Introduction to..."
  - Uneven distribution
    - Many words start with E
    - Few words start with X

# Real Life Hashing IV

- Solutions to issues above?
  - Ignore "The...", "Introduction..."
    - e.g. bucket "The Invisibles" under I
  - Break up E's bucket: 'Ea-Em', 'En-Ez'
  - Merge X's bucket with W/Y
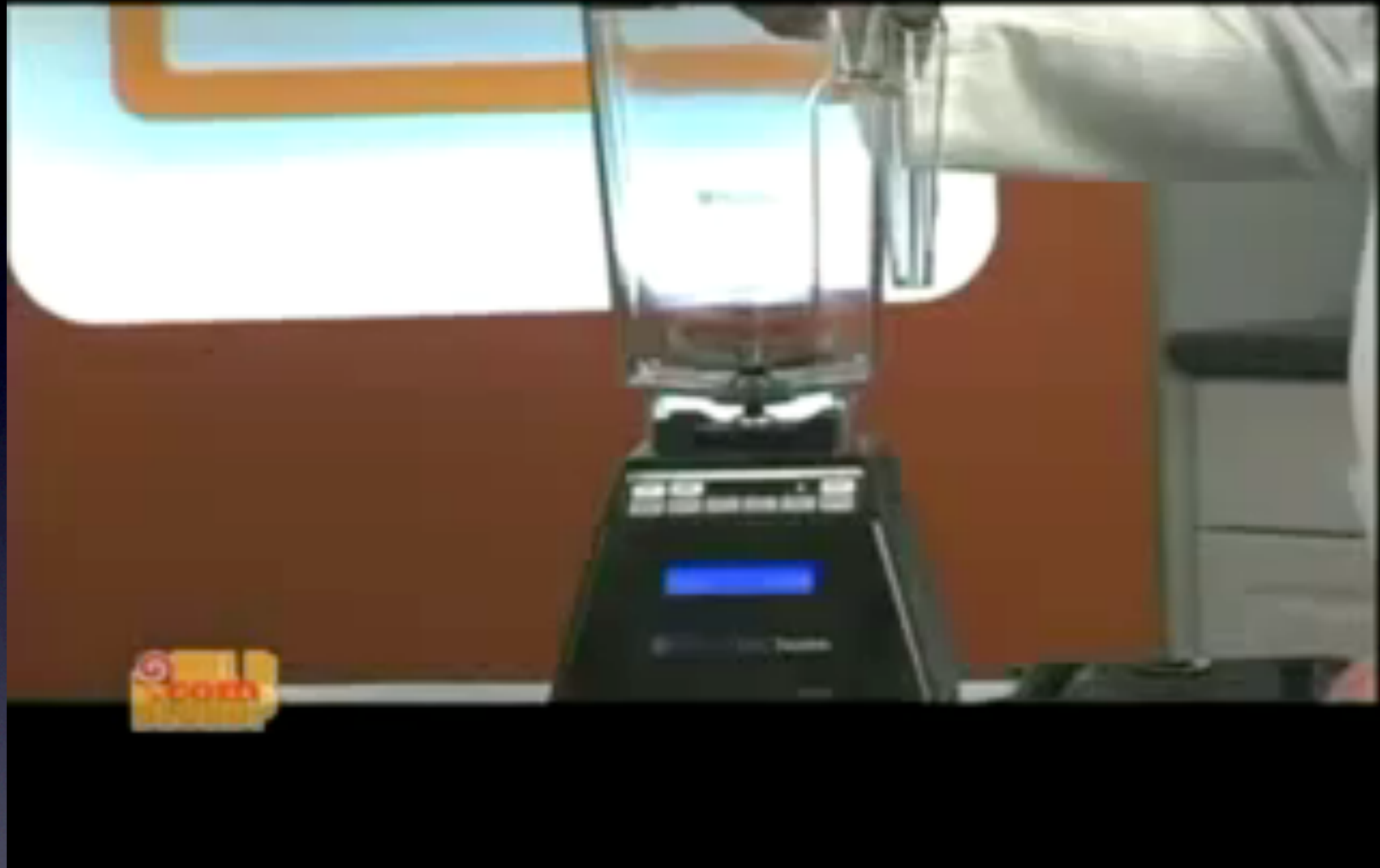- Bucketing function gets hairy :(

# Hashing in Codeworld

- Memory is a block of cells

  - Buckets are numbered 0 to N-1

  - Each bucket is a list of the objects in it

  - Fancy name for the bucketing method:

    **hashing function**

# Hash Functions

- Theory

  - Maps the universe of keys to small (bounded) numbers

- Practice

  - Black magic that allows us to beat the log(N) theoretical bound on a daily basis

# Good Hashing

# Good Hash Functions

- Convenient universe size (16/32/64-bit ints)
- Uniform distribution of keys
  - No obvious bad behavior
- Correct
  - Equal keys always hash to the same value
- Fast

# Hashing Hall of Shame

- String hashing
  - numeric code for first letter
  - sum of numeric code for all letters
    - permutations hash to the same value
  - polynomial value: $\sum str[i] \cdot 256^i$
    - grows without bound

# Hashing Hall of Shame

- String hashing II
  - $(\sum str[i] \cdot 256^i) \bmod 2^{32}$
    - takes $N^2$ to compute
  - $(\sum str[i] \cdot (256^i \bmod 2^{32})) \bmod 2^{32}$
    - only takes first 4 letters into account
    - still sucks for table sizes = powers of 2

# Hashing Wisdom

- Good functions are hard to come up with

- Use built-in functions whenever possible

# Python Hashing 101

- Want hash() to work for your own objects?
  - def __**hash**__(self)
    - hash to a 32-bit number, not -1
- Want your objects as dictionary keys?
  - def __**eq**__(self, other)
    - return True/False (self equals other?)

# Application: Screw Python

- I want lists as dictionary keys!

- Plan:

  1. SuperList object, encapsulating a list

  2. implement __hash__ and __eq__

  3. prepare Turing award acceptance speech

# Behold, it's SuperList!!!

```python
 1 def make32(x):
 2     x = x % (2**32)
 3     if x >= 2**31: x = x - 2**32
 4     return int(x)
 5 class SuperList(object):
 6     def __init__(self, list):
 7         self.list = list
 8     def __hash__(self):
 9         m = 1000003
10         x = 0x345678
11         v = self.list
12         for i in range(len(v)):
13             y = v[i].__hash__()
14             if y == -1: return -1
15             x = make32((x^y)*m)
16             m = make32(m + 82520 + 2*((len(v)-i-1)))
17         x = make32(x+97531)
18         if x == -1:
19             x = -2
20         return x
21     def __eq__(self, other):
22         return self.list.__eq__(other.list)
```

# OMG!! I'm teh one!!!

```
 1 >>> from super_list import SuperList
 2 >>>
 3 >>> k1 = SuperList([1, 2, 3])
 4 >>> k2 = SuperList([1, 2, 3])
 5 >>> k3 = SuperList([4, 5, 6])
 6 >>>
 7 >>> k1 == k2
 8 True
 9 >>> k1 == k3
10 False
11 >>> d = {}
12 >>> d[k1] = 'a'
13 >>> d[k2] = 'b'
14 >>> d[k3] = 'c'
15 >>> print d
16 {<super_list.SuperList object at 0x69870>: 'c',
    <super_list.SuperList object at 0x69930>: 'b'}
17 >>>
18 >>> print d[k1], d[k2], d[k3]
19 b b c
```

# Except not (WTF?!)

```
 1 >>> k1.list.append(4)
 2 >>> k1 == k2
 3 False
 4 >>> k1 == k3
 5 False
 6 >>> hash(k1)
 7 89902565
 8 >>> hash(k3)
 9 448334556
10 >>> d[k1]
11 Traceback (most recent call last):
12   File "<stdin>", line 1, in <module>
13 KeyError: <super_list.SuperList object at 0x69930>
14 >>> d[k2]
15 Traceback (most recent call last):
16    File "<stdin>", line 1, in <module>
17 KeyError: <super_list.SuperList object at 0x698b0>
18 >>> d[k3]
19 'c'
```

# What have we learned?

- Dictionary keys must be immutable

# Hashing Basic Data

- Examine Python's hashing functions for the built-in data types

- Examples of reasonable hash functions, avoiding common pitfalls

- Know your language
  - Especially its cost model

# PyHash: the Plan

```python
 1  def hash(v):
 2      """
 3      A Python implementation that is identical
 4      to the underlying builtin Python function 'hash'
 5      for integers, longs, strings, instances, and tuples thereof.
 6      This returns -1 only when the object is unhashable.
 7      (Floats not yet implemented.)
 8      """
 9      if type(v) == type(1):     return int_hash(v)
10      if type(v) == type(1L):    return long_hash(v)
11      if type(v) == type(" "):   return string_hash(v)
12      if type(v) == type((1,)):  return tuple_hash(v)
13      x = dummy
14      if type(v) == type(x):     return id(v)
15      return -1
```

# PyHash: Short Integers

```python
 1 def make32(x):
 2     """
 3     Convert x into a 32-bit signed integer.
 4     """
 5     x = x % (2**32)
 6     if x >= 2**31:
 7         x = x - 2**32
 8     x = int(x)
 9     return x
10
11 def int_hash(v):
12     if v == -1: v = -2
13     return v
```

# PyHash: Strings

```python
 1 def string_hash(v):
 2     if v == "":
 3         return 0
 4     else:
 5         x = ord(v[0])<<7
 6         m = 1000003
 7         for c in v:
 8             x = make32((x*m)^ord(c))
 9         x ^= len(v)
10         if x == -1:
11             x = -2
12         return x
```

# PyHash: Tuples

```python
 1 def tuple_hash(v):
 2     """
 3   The addend 82520, was selected from the range(0, 1000000) for
 4   generating the greatest number of prime multipliers for tuples
 5   upto length eight:
 6     1082527, 1165049, 1082531, 1165057, 1247581, 1330103, 1082533,
 7     1330111, 1412633, 1165069, 1247599, 1495177, 1577699
 8     """
 9   m = 1000003
10   x = 0x345678
11   for i in range(len(v)):
12       y = v[i].__hash__()  # Invoke built-in python hash
13       if y == -1: return -1
14       x = make32((x^y)*m)
15       m = make32(m + 82520 + 2*((len(v)-i-1)))
16   x = make32(x+97531)
17   if x == -1:
18       x = -2
19   return x
```

# PyHash: Long Integers

```python
 1 def long_hash(v):
 2     sign = 1
 3     if v<0:
 4         v,sign = abs(v),-1
 5     SHIFT = 15                              # for a 32-bit machine
 6     LONG_BIT_SHIFT = 32 - SHIFT
 7     BASE = 1 << SHIFT
 8     MASK = (BASE - 1)
 9     digits = []
10     while v>0:
11         digits.append(v % BASE)
12         v = v>>SHIFT
13     digits.reverse()   # process digits high-order to low-order
14     x = 0
15     for digit in digits:
16         x = (((x << SHIFT) & ~MASK) | ((x >> LONG_BIT_SHIFT) & MASK))
17         x += digit
18         x = make32(x)
19     x = x * sign
20     if x == -1:
21         x = -2
22     return x
```

# And we're done!

- costan@mit.edu

- (617) 230-9694, no voicemail

- AIM: victorcostan

- Google Talk: costan@gmail.com

- 32G-8th Floor

# v. Next

- Notes on contact for \_\_eq\_\_ in Python

  - must check for object type like Java?

- Library example not clear enough, not enough practice with hash functions

- No point in showing all hashing functions... incomprehensible and not good discussion