# CSCI 136
# Data Structures &
# Advanced Programming

Hashing: Loose Ends

# Video Outline

- Growing hashtables
- Choosing an appropriate hashtable size
- Ideal hash function properties and examples
- Revisiting hashtable performance

# Hashtable Size

- Vectors are useful because, when a Vector "runs out of space", the Vector grows
  - It's very clear when we need to grow a vector: excess capacity = 0
- What does it mean for a hashtable to "run out of space"?

- Even ignoring correctness, performance is slowed by "full" hashtables

# Hashtable Size

- The right size for our hashtable will make a trade-off between space and performance
  - We want our table size to be large to minimize collisions (and run/chain lengths) :  ↑  ↓
  - We want our table size to be small to minimize wasted space (empty slots): ↑  ↓
- In addition, we would like some flexibility in case we make a bad initial guess for our size

# Hashtable Fullness: Load Factor

- Suppose a hashtable with M slots stores N elements

- Load factor is a measure of how full the hash table is
    - LF = (# elements) / (# slots) = N / M

- A smaller load factor means the hashtable is less full, which likely gives better performance

# Calculating Load Factor

- To track a hashtable's load factor, we can keep a running count of its elements
  - Every successful `remove()` decrements the count
    - Careful with reserved slots!
    - May want to use chaining if you anticipate many deletes
  - *Some* `put()` operations increment the count
    - Only increment when putting new keys: replacing the value associated with an existing key doesn't change the hashtable's count

- Load factor is then (`count / table.length`)

# Using Load Factor

- Given a hashtable's load factor, what should we do?
  - If the load factor is low, nothing!
    - A low load factor should give good performance
  - If the load factor is high (.6?) , grow our table
    - Increase the number of slots without changing the number of elements (LF = N / M)
- How to grow?
  - Vectors: `ensureCapacity()`
    - Allocate new `Object` array, then copy elements to same index within new (larger) array
      - Does this work for hashtables?

# Load factors

- Idea: always keep load factor below a certain fraction
- Usually .5 to .75
  - Java HashMap (uses external chaining) uses .75
  - structure5 Hashtable (uses linear probing) uses .6

- Plan: keep track of the load factor.  Once it gets too high, make more slots

# Doubling Array

- Cannot just copy values
  - Why?
    - Canonical slot may change
  - Example: suppose (`key.hashCode() == 11`)
    - 11 % 8 = 3;
    - 11 % 16 = 11;
- **Result**: to grow our array, we must recompute the hashcode for each item, then reinsert each item into new array

# Array sizes

- Some people like using hash tables whose size is a prime

- Reason: taking mod the hash table size (if it's a prime) can help "spread out" our items

- Downside: need to find a prime size when "doubling"

# Good Hashing Functions

- Important point:
  - All of our performance hinges on using "good" hash functions that spread keys "evenly"
- Good hash functions:
  - Are fast to compute
  - Uniformly distribute keys across the range
- General rules of thumb?
  - Not really. We almost always have to test "goodness" empirically.

# Example Hash Functions


ASCII TABLE

- What are some feasible hash functions for Strings?
  - Use the first char's ASCII value?
    - 0-255 only
    - Not uniform (some letters more popular than others)
  - Sum of all characters' ASCII values?
    - Not uniform - lots of small words
    - Doesn't give coverage over large array sizes
    - Not good at avoiding collisions – e.g., smile, limes, miles, and slime are all the same
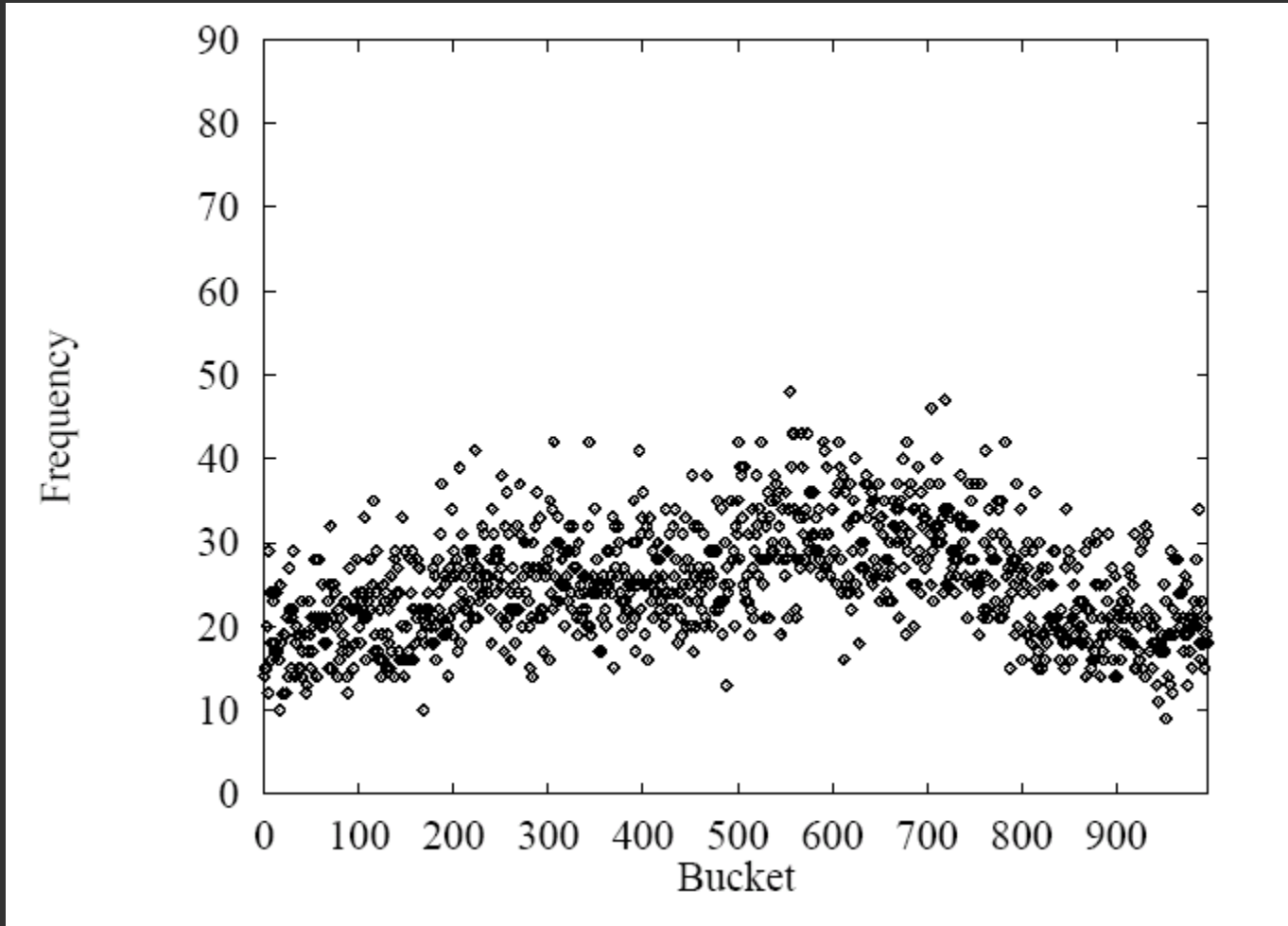    - Let's look at how this works in practice

$$\sum_{i=0}^{s.length()} s.charAt(i)$$

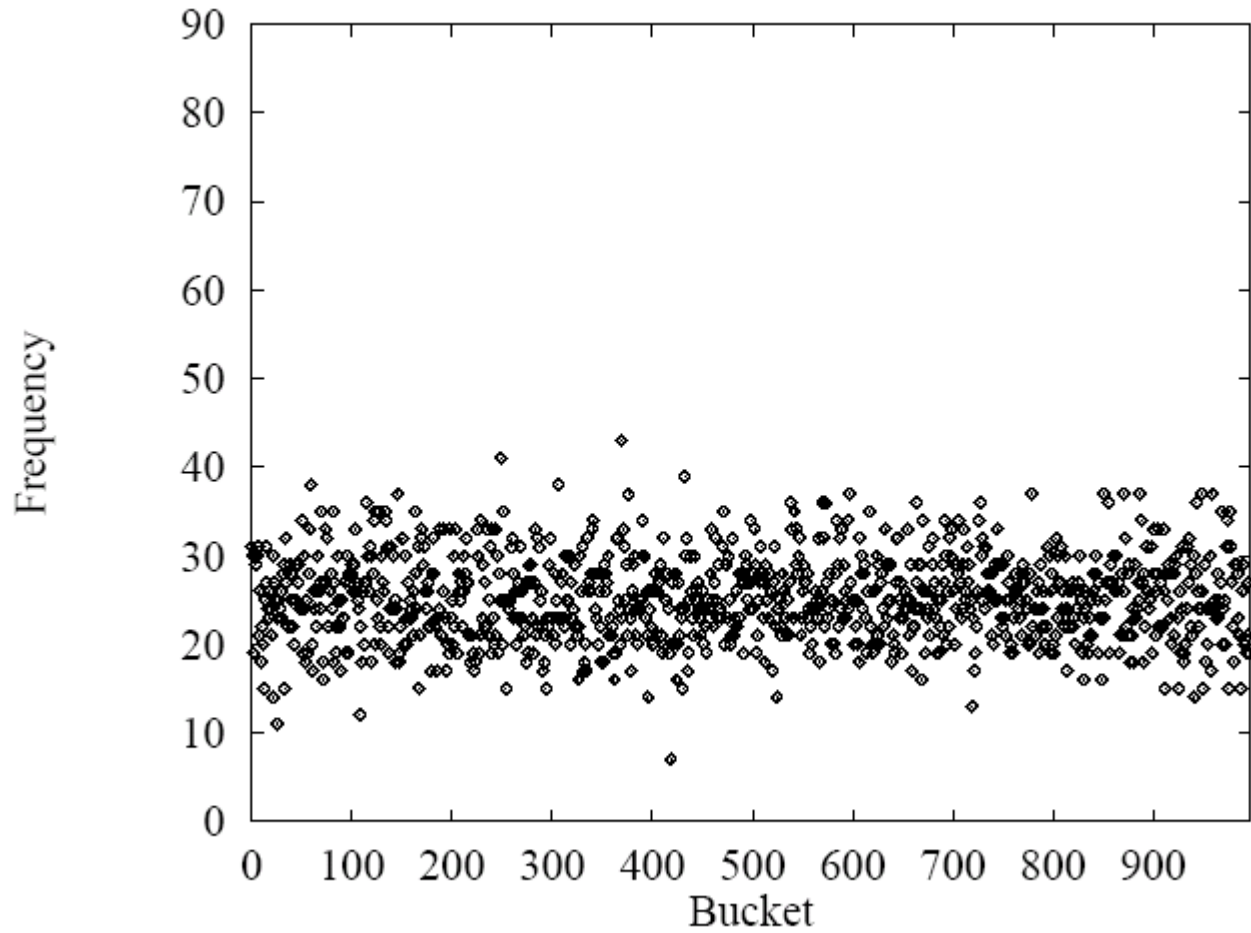Hash of all words in UNIX spelling dictionary (997 buckets)

$$\sum_{i=0}^{s.length()} s.charAt(i) * 2^i$$

Better, but buckets are still pretty uneven (middle quite a bit bigger than ends)

$$\sum_{i=0}^{s.length()} s.charAt(i) * 256^i$$
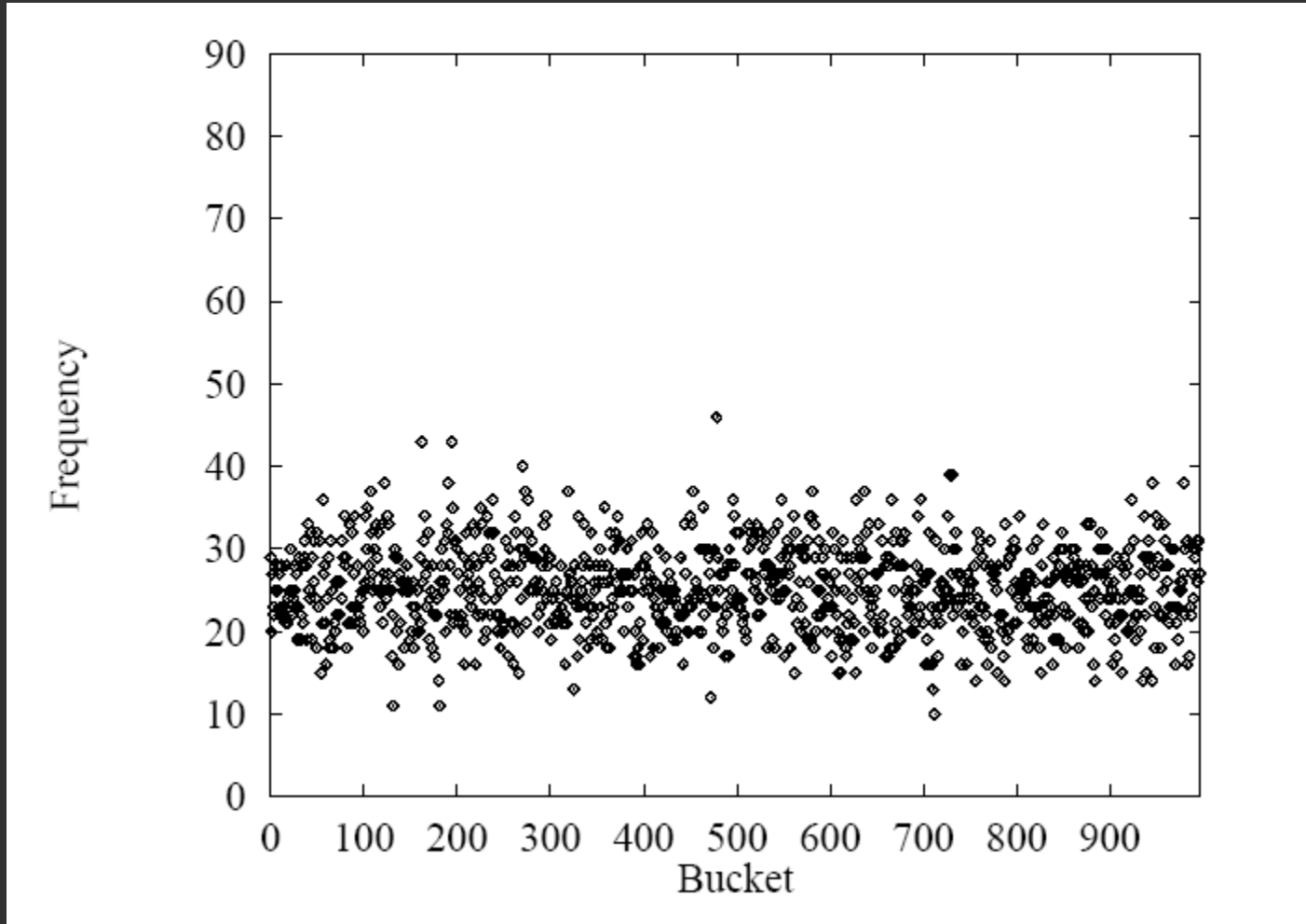
This looks pretty good, but $256^i$ is big…

$$\sum_{i=0}^{s.length()} s.charAt(i) * 31^i$$

$$\sum_{i=0}^{n} s.charAt(i) * 31^{n-i-1}$$

# Other Objects?

- Integer: i.hashCode() is i
  - That might be very bad depending on your data!
  - May want to use another hashCode() method in that case
    - Perhaps a wrapper class to give a new method
- Character, Long similar
- For your own classes: write your own methods!
  - Test empirically to make sure elements are spread out

# Hashtables: O(1) operations?

- How long does it take to compute a String's hashCode?
  - O(s.length())
  - (Doesn't depend on table size)
- Given an object's hash code, how long does it take to find that object?
  - O(run length) or O(chain length) times cost of .equals() method to compare keys

# Impact on performance

- Let's say we have constant load factor
  - Number of slots is a constant factor greater than the number of elements
- And we have a good hash function
  - Spreads objects out "like random"

- Then: an average bucket has constant chain length
- An average bucket is in a run of constant length
- (Worst case is O(log n) for both---but this is very rare)
- Usually we say that hash tables have O(1) performance

# Summary

|  | put | get | space |
|---|---|---|---|
| unsorted vector | O(n) | O(n) | O(n) |
| unsorted list | O(n) | O(n) | O(n) |
| sorted vector | O(n) | O(log n) | O(n) |
| balanced BST | O(log n) | O(log n) | O(n) |
| hashtable | O(1)* | O(1)* | O(n)* |