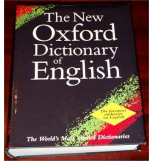


Lecture 15

Dictionaries and Hash Tables



EECS 281: Data Structures & Algorithms

Dictionary ADT

A container of items (key/value pairs) that supports two basic operations

- **Insert** a new item
- **Search (retrieve)** an item with a given key

Two primary uses

- A set of things: check if something is in the set
- Key-Value storage: look up values by keys

4

Dictionary ADT Operations

- Desirable Operations
 - **Insert** a new item
 - **Search** for item(s) having a given key
 - **Remove** a specified item
 - **Sort** the dictionary
 - **Select** the k^{th} largest item in a dictionary
 - **Join** two dictionaries
- Other basic container operations: construct, test if empty, destroy, copy...

5

Containers with key lookup

Identifying a container with fast search and fast insert for arbitrary <key, value> pairs

- Sorted Vectors: Insert will be $O(n)$
- Unsorted Vectors: Search will be $O(n)$
- Sorted/Unsorted Linked Lists: Search is $O(n)$
- Binary Search Tree (std::map<>)
 - Search: average $O(\log n)$, worst-case $O(n)$
 - Insert: average $O(\log n)$, worst-case $O(n)$
- Hash Table (std::unordered_map<>)
 - Search: average $O(1)$, worst case $O(n)$
 - Insert: average $O(1)$, worst case $O(n)$

6

What if the set of keys is small?

- Example: calendar for 1..n days
 - n could be 365 or 366
 - Can look up a particular day in $O(1)$ time
 - Every day is represented by a bucket, i.e., some container
- If we have a range of integers that fits into memory, everything is easy
 - **What if we don't?**

7

Hashing Defined

Locate items in a table by key

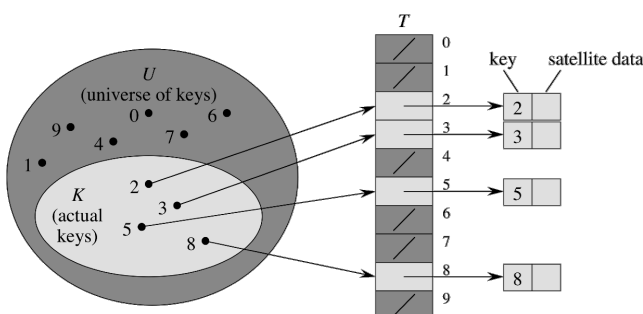
- Use arithmetic operations to calculate a table index (bucket) from a given key

Need

- **Translation:** converts a search key into an integer
- **Compression:** limits an integer to a valid index
- **Collision resolution:** resolves search keys that hash to same table index

10

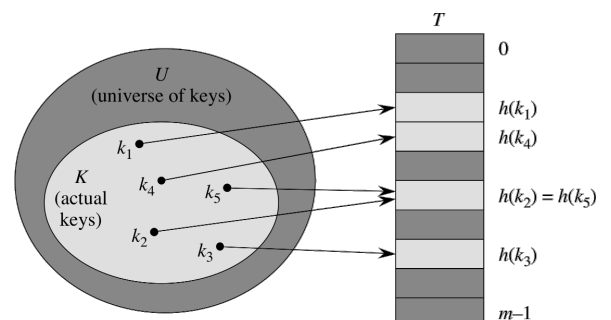
Direct Addressing



Each key maps to an element of the array

11

Hashed Addressing



Only actual keys map to an element of the array

12

Hash Function

Translation: $t(key) \Rightarrow hashint$

- Converts the key into an integer

Compression: $c(hashint) \Rightarrow \text{table index}$

- Maps the hashed integer into the range $[0, M)$

A hash function combines translation and compression:

$$h(key) \Rightarrow c(t(key)) \Rightarrow \text{table index}$$

Note: `std::hash<>` provides only translation, not compression

13

Translating Floating Points

- key between 0 and 1: $[0, 1)$

$$h(key) = \lfloor key * M \rfloor$$

- key between s and t : $[s, t)$

$$h(key) = \lfloor (key - s) / (t - s) * M \rfloor$$

- Try: range = $[1.38, 6.75)$, $M = 13$
compute $h(3.65)$

15

Translating Strings

- Simple hash sums ASCII character codes
- Problem Case: stop, tops, pots, spot
 - Sum of each is equivalent
 - All will map to same hash table address
 - Multiple collisions
- Solution: Character position is important
 - Consider decimal numbers

$$123 = 1 * 10^2 + 2 * 10^1 + 3 * 10^0$$

$$321 = 3 * 10^2 + 2 * 10^1 + 1 * 10^0$$

18

Compression

$c(hashint) \Rightarrow \text{index in range } [0, M)$

- $hashint$ may be < 0 or $\geq M$

Division Method

$\lfloor hashint \rfloor \bmod M$, where M is prime

MAD (multiply and divide) Method

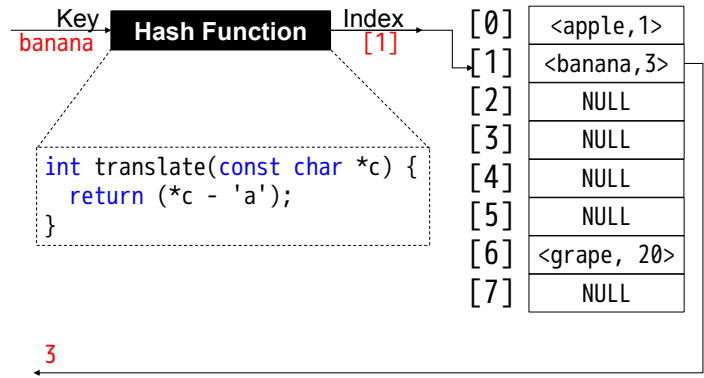
$\lfloor a * hashint + b \rfloor \bmod M$, where a and b are prime

Use this method when you can't control M

Note: $a \bmod M$ must not equal 0!

20

Hashing Example



14

Translating Integers

Modular hash function

$$t(key) = key$$

$$h(key) = c(t(key)) = key \bmod M$$

- Great if keys randomly distributed
 - Often, keys are not randomly distributed
 - Example: midterm bubbles scores all multiples of 2.5
 - Example: pick a number from 1 to 10
- Don't want to pick a bad M
 - BAD: M and key have common factors

17

Better Translation: Rabin Fingerprint

- Instead of adding up character codes, view strings as *decimal numbers*
 - Characters: 'T', 'O', 'M', ' ', 'M', ...
 - ASCII codes: 84, 79, 77, 32, 77, ...
 - Running fingerprints:

T	84
TO	$10 * 84 + 79$
TOM	$10 * (10 * 84 + 79) + 77$
TOM	$10 * (10 * (10 * 84 + 79) + 77) + 32, \dots$
 - Base 10 is used for illustration only (use larger numbers)
- **Shuffling the chars usually changes result**

19

Hash Function

Translation: $t(key) \Rightarrow hashint$

- Converts the key into an integer

Compression: $c(hashint) \Rightarrow \text{table index}$

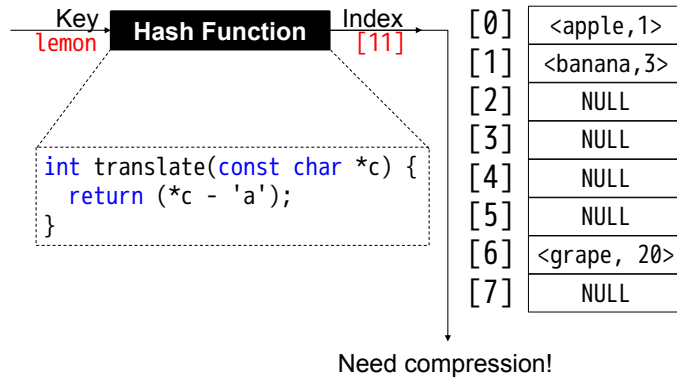
- Maps the hashed integer into the range $[0, M)$

A hash function combines translation and compression:

$$h(key) \Rightarrow c(t(key)) \Rightarrow \text{table index}$$

21

Hashing Example



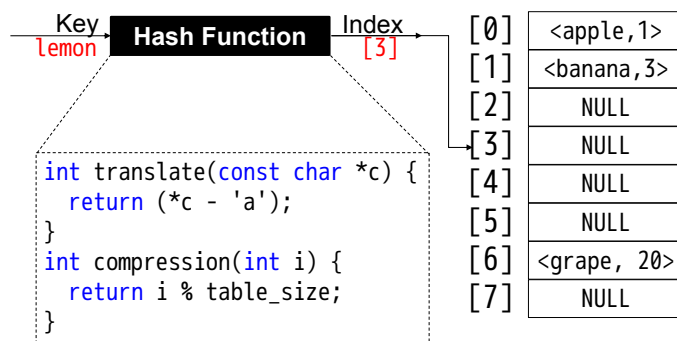
22

Hash Table Size

- Table of size M
 - About the number of elements expected
 - If unsure, guess high
- Hash function must return keys as integers in range $[0, M)$

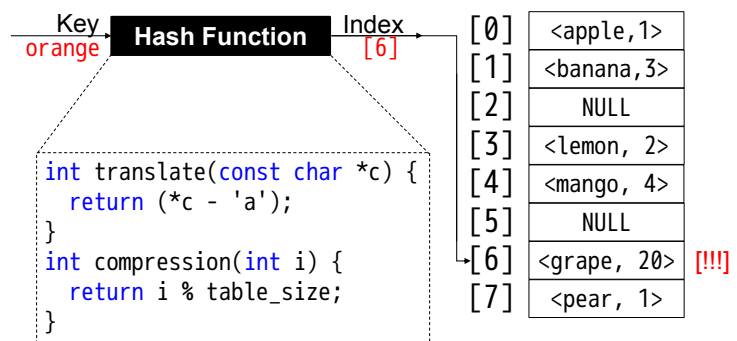
23

Index Compression



24

Collision



25

Hash Tables Summarized

- Efficient ADT for insert, search, and remove
- Hashing a key $h(key) \Rightarrow \text{table index}$
 - Maps key to table index in two steps
 - Translates key into an integer
 $t(key) \Rightarrow \text{hashint}$
 - Compression maps integer into range $[0, M)$
 $c(\text{hashint}) \Rightarrow \text{table index}$
- Therefore, $h(key) \equiv c(t(key))$

26

Good Hash Functions

- Benefits of hash tables depend on having “good” hash functions
- Must be easy to compute
 - Will compute a hash for every key
 - Will compute same hash for same key
- Should distribute keys evenly in table
 - Will minimize *collisions*
 - Collision: two keys map to same index
- Trivial, poor hash function: $h(key) \{ \text{return } 0; \}$
 - Easy to compute, but poor distribution maximizes collisions

29

Complexity of Hashing

For simplicity, assume perfect hashing (no collisions)

- What is cost of **insertion**? $O(1)$
- What is cost of **search**? $O(1)$
- What is cost of **removal**? $O(1)$

Wouldn't it be nice to live in a perfect world?

- Recall the Pigeonhole Principle

30

Real-World Hash Tables

- C++ hash table containers (STL, C++11+)
 - `unordered_set<>`, `unordered_multiset<>`
 - `unordered_map<>`, `unordered_multimap<>`
- Database indices are built on hash tables
- Compilers use a hash table for identifiers

31

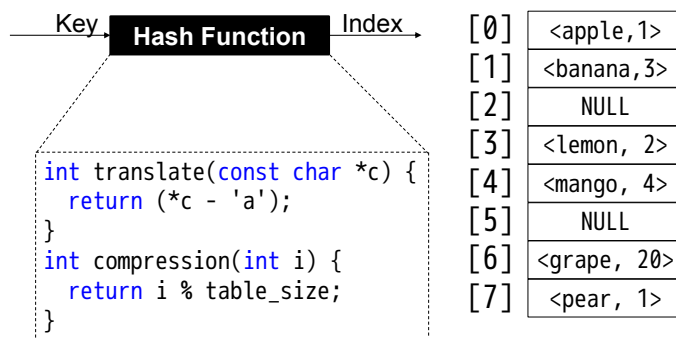
Major Uses of Hash Tables

- Sets of ints, strings, images, class objects
 - Check set membership
 - Detect/filter duplicates
 - Find unique elements
 - Find matching elements, e.g., $a + b = 0$
- Key-value storage: look up values by keys
 - Count things: `value_type = int`
 - Maintain *sparse* vectors: `key_type = int`
 - Maintain linked structures: `key_type = value_type`



34

Sorting a Hash Table?



35

Dictionary ADTs & Hashing

Hashing is an **efficient** implementation of:

- Insert* a new item
- Search* for an item (or items) having a given key
- Remove* a specified item

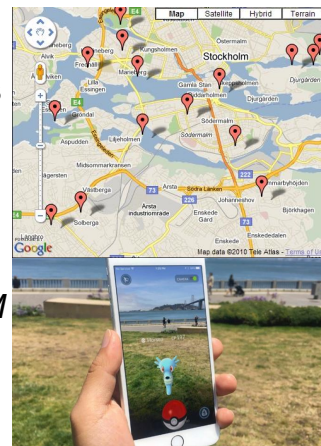
Hashing is an **inefficient** implementation of:

- Select* the k^{th} largest item in a dictionary
- Sort* the items in the dictionary

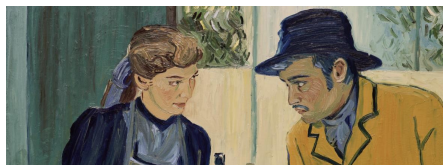
36

Composite Hash Functions

- Example: hang some info over geo-locations
 - Hash {long, lat} pairs
- How do we combine two hash functions?**
- $H(\{x, y\}) = (h(x) + p * h(y)) \% M$
- How do we hash class objects?



38



Is $(h_1 + p * h_2) \% M$ a good hash combiner?

- Recall properties of good hash functions
 - Fast computation
 - Even distribution of values
 - No easily-predictable collisions**

$$((h_1 + p * h_2) + p * h_3) \% M = ((h_1 + p * h_3) + p * h_2) \% M$$

- A better combiner for n values ($i = 0 \dots n - 1$)

39

When Not To Use Hash Tables

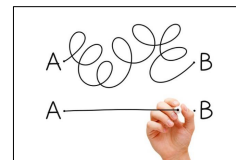
- When keys are small ints, use **bucket arrays**
 - Some keys can be coerced to small integers: `enums`, `simple fractions`, `months`
- For static data, set membership, or lookup
 - Consider sorting + binary search
- For key-value storage where traversals are needed but not lookup (e.g., sparse vec)
 - Store key-value pairs in a list (or vector)



41

When Not To Use Hash Tables

- In many applications, it is tempting to use `unordered_set<>` or `unordered_map<>`, but
 - Significant space overhead of nested containers
 - Every access computes hash function
 - $O(n)$ worst-case time (STL implementations)



40

Word Count Demo

From a web browser:
bit.ly/eecs281-wordcount-demo

From a terminal:
`wget bit.ly/eecs281-wordcount-demo -O wordcount.cpp`

At the command line:
`g++ -std=c++1z wordcount.cpp -o wordcount`
`./wordcount`

Enter filename: wordcount.cpp



41