# Algorithms
## (Hash Tables)

**Pramod Ganapathi**
Department of Computer Science
State University of New York at Stony Brook

March 19, 2021

## Contents

- (Sorted) Sets, Maps, Multisets, Multimaps
- Hash Tables
- Hash Functions
- Collisions
- Collision-Avoiding Techniques
  - Separate Chaining
  - Open Addressing

## Dictionary ADTs

These ADTs are collection of items, where,
each item can be a $key$ or a $(key, value)$ pair.

| ADT | Item | Ordered? | Duplicates? | Implementation |
|------|------|----------|-------------|----------------|
| Set | $key$ | ✗ | ✗ | Hash table |
| Sorted set | $key$ | ✓ | ✗ | Balanced tree |
| Multiset | $key$ | ✗ | ✓ | Hash table |
| Sorted multiset | $key$ | ✓ | ✓ | Balanced tree |
| Map | $(key, value)$ | ✗ | ✗ | Hash table |
| Sorted map | $(key, value)$ | ✓ | ✗ | Balanced tree |
| Multimap | $(key, value)$ | ✗ | ✓ | Hash table |
| Sorted multimap | $(key, value)$ | ✓ | ✓ | Balanced tree |

# Set ADT (java.util.Set interface)

| Method | Functionality |
|---|---|
| add(e) | Adds the element $e$ to $S$ (if not already present). |
| remove(e) | Removes the element $e$ from $S$ (if it is present). |
| contains(e) | Returns whether $e$ is an element of $S$. |
| iterator() | Returns an iterator of the elements of $S$. |
| addAll(T) | Updates $S$ to also include all elements of set $T$, effectively replacing $S$ by $S \cup T$. |
| retainAll(T) | Updates $S$ so that it only keeps those elements that are also elements of set $T$, effectively replacing $S$ by $S \cap T$. |
| removeAll(T) | Updates $S$ by removing any of its elements that also occur in set $T$, effectively replacing $S$ by $S - T$. |

- Set = unordered set; Map = unordered map.
  java.util.HashSet is an implementation of the set ADT.
  java.util.HashMap is an implementation of the map ADT.

# Sorted set ADT (java.util.SortedSet interface)

| Method | Functionality |
|---|---|
| `first()` | Returns the smallest element in $S$. |
| `last()` | Returns the largest element in $S$. |
| `ceiling(e)` | Returns the smallest element $\geq e$. |
| `floor(e)` | Returns the largest element $\leq e$. |
| `lower(e)` | Returns the largest element $< e$. |
| `higher(e)` | Returns the smallest element $> e$. |
| `subSet(e1,e2)` | Returns an iteration of all elements greater than or equal to $e1$, but strictly less than $e2$. |
| `pollFirst()` | Returns and removes the smallest element in $S$. |
| `pollLast()` | Returns and removes the largest element in $S$. |

- java.util.TreeSet is an implementation of the sorted set ADT. java.util.TreeMap is an implementation of the sorted map ADT.

# Multiset ADT

| Method | Functionality |
|---|---|
| `add(e)` | Adds a single occurrences of $e$ to the multiset. |
| `contains(e)` | Returns true if the multiset contains an element $= e$. |
| `count(e)` | Returns the number of occurrences of $e$ in the multiset. |
| `remove(e)` | Removes a single occurrence of $e$ from the multiset. |
| `remove(e, n)` | Removes $n$ occurrences of $e$ from the multiset. |
| `size()` | Returns the number of elements of the multiset (including duplicates). |
| `iterator()` | Returns an iteration of all elements of the multiset (repeating those with multiplicity greater than one). |

- Java does not include any form of a multiset.
  Guava = Google Core Libraries for Java.
  Guava's Multiset is an implementation of the multiset ADT.
  Guava's Multimap is an implementation of the multimap ADT.
- Similarly, one can define sorted multiset ADT

# Hash Tables

# Hash tables

- A hash table is an efficient dictionary data structure to implement a set/multiset/map/multimap.
- A hash table performs put, remove, and get operations in constant expected time.
- Hashing is the implementation of hash tables.

# Balanced search trees vs. Hash tables

- Balanced search tree $\Leftrightarrow$ sorted, Hash table $\Leftrightarrow$ unsorted
- Worst $=$ worst-case, avg. $=$ expected time (useful in practice)

| Operations | | Balanced tree (worst) | Hash table (avg.) | Hash table (worst) |
|---|---|---|---|---|
| Sorting-unrelated operations | Insert | $\mathcal{O}(\log n)$ | $\mathcal{O}(1)$ | $\mathcal{O}(n)$ |
| | Delete | $\mathcal{O}(\log n)$ | $\mathcal{O}(1)$ | $\mathcal{O}(n)$ |
| | Search | $\mathcal{O}(\log n)$ | $\mathcal{O}(1)$ | $\mathcal{O}(n)$ |
| Sorting-related operations | Sort | $\mathcal{O}(n)$ | ✗ | |
| | Minimum | $\mathcal{O}(\log n)$ | ✗ | |
| | Maximum | $\mathcal{O}(\log n)$ | ✗ | |
| | Predecessor | $\mathcal{O}(\log n)$ | ✗ | |
| | Successor | $\mathcal{O}(\log n)$ | ✗ | |
| | Range-Minimum | $\mathcal{O}(\log n)$ | ✗ | |
| | Range-Maximum | $\mathcal{O}(\log n)$ | ✗ | |
| | Range-Sum | $\mathcal{O}(n)$ | ✗ | |

## Applications of hash tables

- Web page search using URLs
- Password verification
- Symbol tables in compilers
- Filename-filepath linking in operating systems
- Plagiarism detection using Rabin-Karp string matching algorithm
- English dictionary search
- Used as part of the following concepts:
  - finding distinct elements
  - counting frequencies of items
  - finding duplicates
  - message digests
  - commitment
  - Bloom filters

# Map

- A map is a collection of key-value pairs $(k, v)$, where, keys are unique.

| Key | Value |
|---|---|
| User ID | User record |
| Employee ID | Employee record |
| Student ID | Student record |
| Patient ID | Patient record |
| Profile ID | Person details |
| Order ID | Order details |
| Transaction ID | Transaction details |
| URL | Web page |
| Full file name | File |

# Hash tables

- A hash table is an efficient implementation of a set or map, i.e., insert, delete, and search operations take constant expected time.
- Example: Suppose we store (name, favorite color) pairs
  We place the key-value pairs in the cells of the hash table array

# Hash tables

## Questions

- Why do we need to think in terms of (key, value) pairs? Why not $k$-tuples?
- How are keys of arbitrary objects mapped to array indices which are whole numbers?
- A hash table is a data structure of finite size. How can an infinite number of keys be mapped to a finite number of indices?
- Can there be collisions during mapping? That is, isn't there a nonzero chance that different keys get mapped to the same index?
- Is there a relation between the table size $N$ and the number of elements $n$?
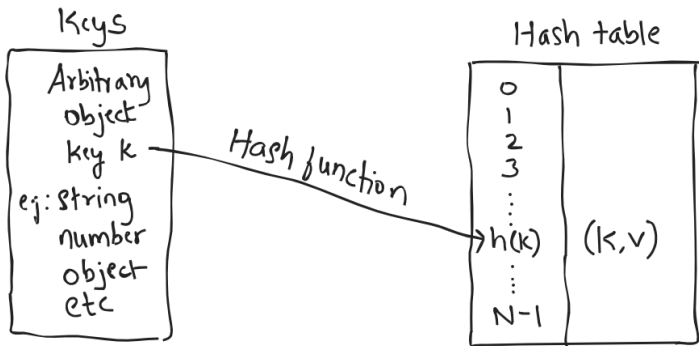
# Hash Functions

# Hash functions

## Questions

1. How are keys of arbitrary objects mapped to array indices which are whole numbers?
2. A hash table is a data structure of finite size. How can an infinite number of keys be mapped to a finite number of indices?
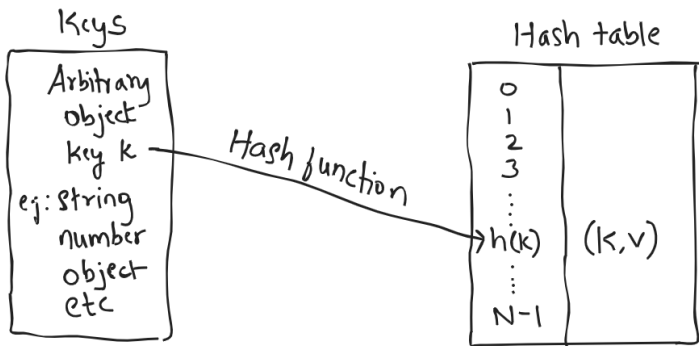
## Solution

- Answer: Hash function

# Hash functions

## Solution (continued)



- A hash function is a mapping from arbitrary objects to the set of indices $[0, N-1]$.
- A hash function stores key-value pair $(k, v)$ in array $A[h(k)]$.
- A hash function is good when it is easy to compute, fast to compute, and leads to few collisions.

# Hash functions

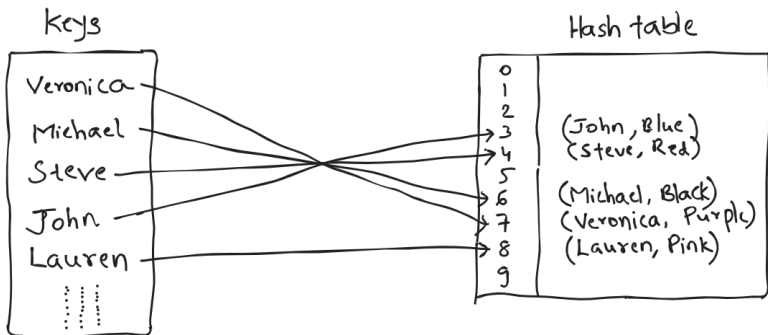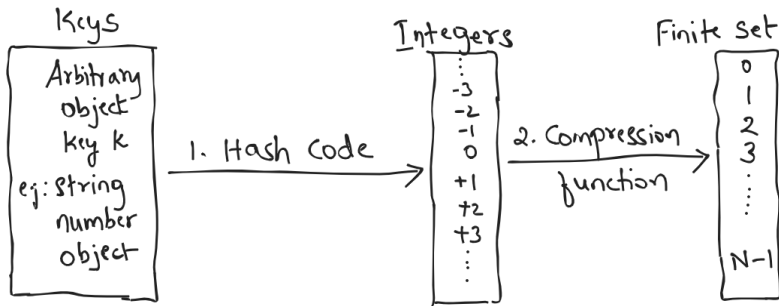## Solution (continued)



- A hash function is a mapping from arbitrary objects to the set of indices $[0, N-1]$.
- A hash function stores key-value pair $(k, v)$ in array $A[h(k)]$.
- A hash function is good when it is easy to compute, fast to compute, and leads to few collisions.

# Hash functions

## Solution (continued)



- For modularity, assume a hash function consists of two stages:
  1. Hash code
  2. Compression function

  Advantage: The hash code portion of the computation is independent of a specific hash table size

# Hash codes

- Consider bits as integer.
  Hashcode(byte | short | char) = 32-bit int $\qquad \triangleright$ upscaling
  Hashcode(float) = 32-bit int $\qquad \triangleright$ change representation
  Hashcode(double) = 32-bit int $\qquad \triangleright$ downscaling
  Hashcode$(x_0, x_1, \ldots, x_{n-1}) = x_0 + x_1 + \cdots + x_{n-1}$ $\qquad \triangleright$ sum
  Hashcode$(x_0, x_1, \ldots, x_{n-1}) = x_0 \oplus x_1 \oplus \cdots \oplus x_{n-1}$ $\qquad \triangleright$ xor
- Polynomial hash codes.
  Hashcode$(x_0, x_1, \ldots, x_{n-1}) =$
  $$x_0 a^{n-1} + x_1 a^{n-2} + \cdots + x_{n-2}a + x_{n-1}$$ $\qquad \triangleright$ polynomial
- Cyclic-shift hash codes.
  Hashcode$_k(x) =$ Rotate$(x, k$ bits$)$ $\qquad \triangleright$ cyclic-shift
  e.g.: Hashcode$_2(111000) = 100011$

# Compression functions

A good compression function minimizes the number of collisions for a given set of distinct hash codes.

- Division method.

  $\boxed{\text{Compression}(i) = i \ \% \ N}$        ▷ remainder

  $N \geq 1$ is the size of the bucket array.
  Often, $N$ being prime "spreads out" the distribution of primes.
  Ex. 1: Insert codes $\{200, 205, \ldots, 600\}$ into $N$-sized array.
  Which is better: $N = 100$ or $N = 101$?
  Ex. 2: Insert multiple codes $\{aN + b\}$ into $N$-sized array.
  Which is better: $N =$ prime or $N =$ non-prime?

- Multiply-Add-and-Divide (MAD) method.

  $\boxed{\text{Compression}(i) = ((ai + b) \ \% \ p) \ \% \ N}$      ▷ remainder

  $N \geq 1$ is the size of the bucket array.
  $p$ is a prime number larger than $N$.
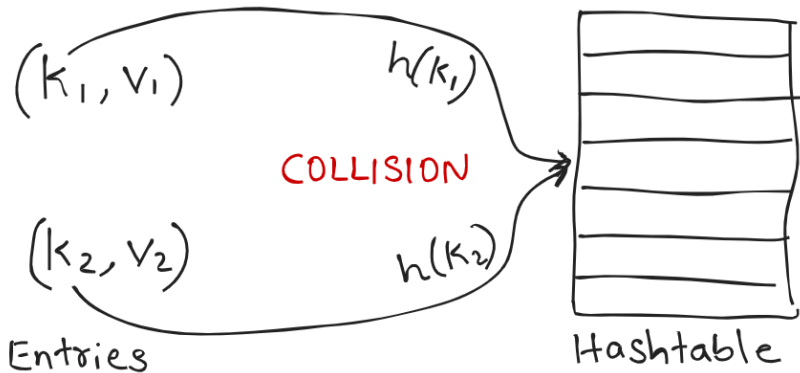  $a, b$ are random integers from the range $[0, p - 1]$ with $a > 0$.
  Usually eliminates repeated patterns in the set of hash codes.

# Collisions

Suppose you want to insert two entries $(k_1, v_1)$ and $(k_2, v_2)$ into a hashtable such that $h(k_1) = h(k_2)$. This is called collision as you cannot insert both the entries at the same location.

So, we need to handle collisions.

## Collision-handling schemes

There are two major collision-handling schemes or collision-resolution strategies.

| Collision-handling scheme | Features |
|---|---|
| Separate chaining | Extra space (for secondary data structures) |
| | Simpler implementation |
| Open addressing | No extra space |
| | More complicated implementation |

# Separate chaining

- Have each bucket $A[j]$ store its own secondary container.
- We use secondary data structures (e.g. array list, linked list, balanced search trees, etc) for each bucket.

# Separate chaining (via arraylist/linkedlist)

---

Put($(key, value)$)

1. $hash \leftarrow$ Hash($key$)
2. $A[hash]$.AddLast($(key, value)$)     ▷ $A[hash]$ is a linked list

Get($key$)

1. $hash \leftarrow$ Hash($key$)
2. **return** $A[hash]$.Get($key$)     ▷ returns value

Remove($key$)

1. $hash \leftarrow$ Hash($key$)
2. **return** $A[hash]$.Remove($key$)     ▷ returns removed value

---

# Open addressing

- All entries are stored in the bucket array itself.
- Strict requirement: Load factor must be at most 1.
- Useful in applications where there are space constraints,
  e.g.: smartphones and other small devices.
- Iteratively search the bucket $A[(\text{HASH}(key) + f(i)) \% N]$
  for $i = 0, 1, 2, 3, \ldots$ until finding an empty bucket.

| Scheme | Function |
|---|---|
| Linear probing | $f(i) = i$ |
| Quadratic probing | $f(i) = i^2$ |
| Double hashing | $f(i) = i \cdot \text{HASH2}(key)$ |
| | e.g. $\text{HASH2}(key) = p - (key \% p)$ for prime $p < N$. |
| | Here, $N$ should be a prime number. |
| Random generator | $f(i) = \text{RANDOM}(i, \text{HASH}(key))$ |

# Linear probing: Put

- Suppose $\text{HASH}(key) = key \% 10$

| Put | | | Array | | | | | | | | | |
| Key | $\rightarrow$ | Hash | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 18 | $\rightarrow$ | 8 | | | | | | | | | 18 | |
| 41 | $\rightarrow$ | 1 | | 41 | | | | | | | 18 | |
| 22 | $\rightarrow$ | 2 | | 41 | 22 | | | | | | 18 | |
| 32 | $\rightarrow$ | 2 | | 41 | 22 | | | | | | 18 | |
| (2 probes) | | | | 41 | 22 | 32 | | | | | 18 | |
| 98 | $\rightarrow$ | 8 | | 41 | 22 | 32 | | | | | 18 | |
| (2 probes) | | | | 41 | 22 | 32 | | | | | 18 | 98 |
| 58 | $\rightarrow$ | 8 | | 41 | 22 | 32 | | | | | 18 | 98 |
| | | | | 41 | 22 | 32 | | | | | 18 | 98 |
| (3 probes) | | | 58 | 41 | 22 | 32 | | | | | 18 | 98 |
| 78 | $\rightarrow$ | 8 | How many probes are required to insert 78? | | | | | | | | | |

# Linear probing: Remove

- Suppose $\text{HASH}(key) = key \% 10$

| Remove Key | Array | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| − | 58 | 41 | 22 | 32 | 78 | 19 | | | 18 | 98 |
| 58 | 58 | 41 | 22 | 32 | 78 | 19 | | | 18 | 98 |
| | | 41 | 22 | 32 | 78 | 19 | | | 18 | 98 |
| 19 | | 41 | 22 | 32 | 78 | 19 | | | 18 | 98 |

<div style="text-align:center; color:red;">Hence, we cannot simply remove a found entry.</div>

| Remove Key | Array | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
| − | 58 | 41 | 22 | 32 | 78 | 19 | | | 18 | 98 |
| 58 | 58 | 41 | 22 | 32 | 78 | 19 | | | 18 | 98 |
| | 58 | 41 | 22 | 32 | 78 | 19 | | | 18 | 98 |
| 19 | 58 | 41 | 22 | 32 | 78 | 19 | | | 18 | 98 |
| | 58 | 41 | 22 | 32 | 78 | 19 | | | 18 | 98 |

<div style="text-align:center; color:red;">Replace the deleted entry with the defunct object.</div>

## Linear probing

---

$\text{Put}((key, value))$

1. $hash \leftarrow \text{Hash}(key);\ i \leftarrow 0$
2. **while** $(hash + i)\ \%\ N \neq null$ **and** $i < N$ **do** $i \leftarrow i + 1$
3. **if** $i = N$ **then throw** Bucket array is full
4. **else** $A[(hash + i)\ \%\ N] \leftarrow (key, value)$

---

$\text{Get}(key)$

1. $hash \leftarrow \text{Hash}(key);\ i \leftarrow 0$
2. **while** $(hash + i)\ \%\ N \neq null$ **and** $i < N$ **do**
3.    $index \leftarrow (hash + i)\ \%\ N$
4.    **if** $A[index].key = key$ **then return** $A[index].value$
5.    $i \leftarrow i + 1$
6. **return** $null$

---

$\text{Remove}(key)$

1. $index \leftarrow \text{FindSlotForRemoval}(key)$
2. **if** $index < 0$ **then return** $null$
3. $value \leftarrow A[index].value;\ A[index] \leftarrow defunct;\ n \leftarrow n - 1$
4. **return** $value$

---

# Separate chaining, Open addressing: Complexity

- Suppose $N =$ bucket array size and $n =$ number of entries.
- Ratio $\lambda = n/N$ is called the load factor of the hash table.
- If $\lambda > 1$, rehash. Make sure $\lambda < 1$.
- Assuming good hash function, expected size of bucket is $\mathcal{O}\left(\lceil \lambda \rceil\right)$.
- Separate chaining: Maintain $\lambda < 0.75$
  Open addressing: Maintain $\lambda < 0.5$
- Assuming good hash function and $\lambda \in \mathcal{O}\left(1\right)$,
  complexity of put, get, and remove is $\mathcal{O}\left(1\right)$ expected time.