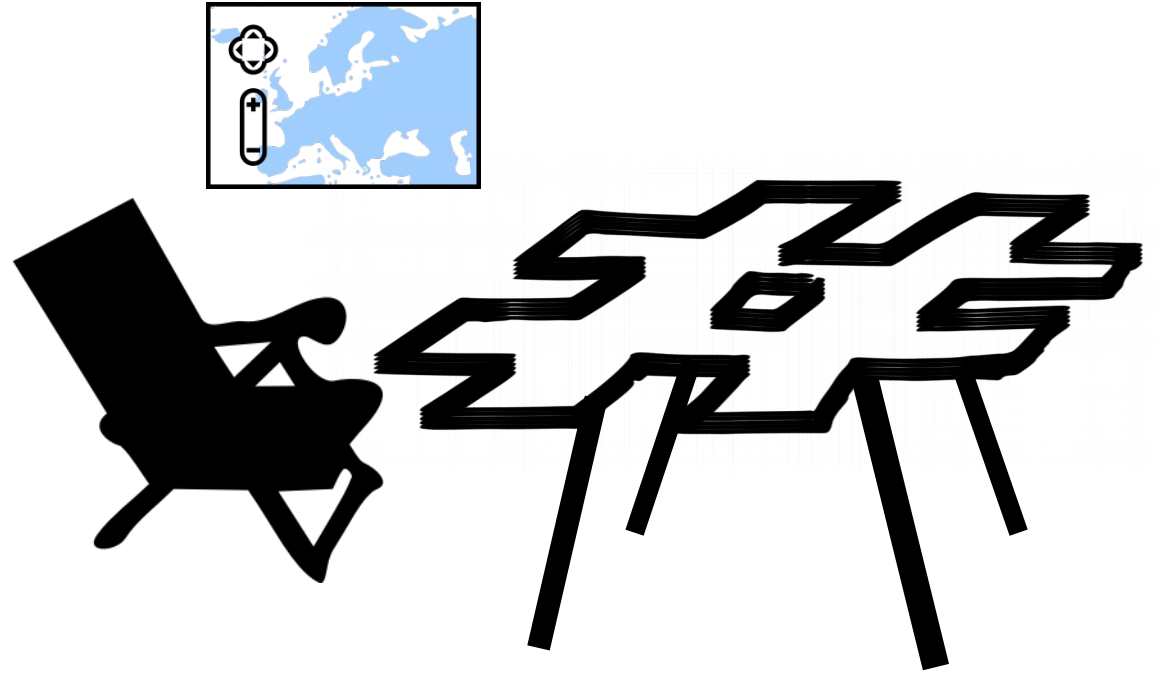


CS 240
Fall 2015

Mike Lam, Professor



Maps and Hash Tables

ADTs

- List
- Set
- Stack
- Queue
 - Deque
 - Priority Queue

New ADT: Map

- Map: key → value
- Unique keys, non-unique values
- Other names:
 - Dictionary
 - Associative array
- Many applications
 - Database: student ID# → student info
 - DNS: domain name → IP address
 - OS: process ID → process record
 - Namespace: variable → value

Map ADT

- **set**(M, k, v) modify value for key
- **get**(M, k) retrieve value for key
- **contains**(M, k) return true if key has a mapping
- **remove**(M, k) remove key from map
- **size**(M) return # of keys
- **clear**(M) remove all mappings
- **key_set**(M) return a set of all keys
- **value_set**(M) return a set of all values

Map Implementations

- Unsorted array
 - Insert: Lookup: Modify:
- Sorted array
 - Insert: Lookup: Modify:
- Skip list / AVL tree
 - Insert: Lookup: Modify:

Map Implementations

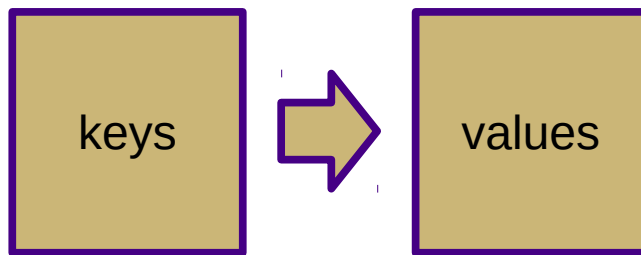
- Unsorted array
 - Insert: $O(n)$ Lookup: $O(n)$ Modify: $O(n)$
- Sorted array
 - Insert: $O(n)$ Lookup: $O(\log n)$ Modify: $O(\log n)$
- Skip list / AVL tree
 - Insert: $O(\log n)$ Lookup: $O(\log n)$ Modify: $O(\log n)$

Map Implementations

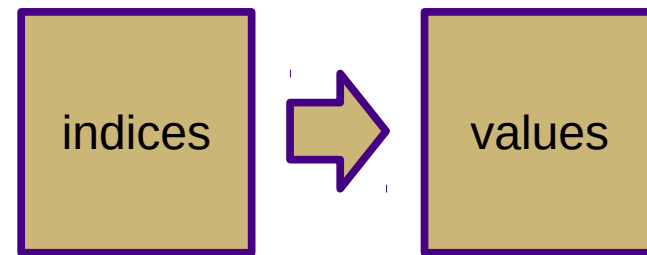
- Unsorted array
 - Insert: $O(n)$ Lookup: $O(n)$ Modify: $O(n)$
- Sorted array
 - Insert: $O(n)$ Lookup: $O(\log n)$ Modify: $O(\log n)$
- Skip list / AVL tree
 - Insert: $O(\log n)$ Lookup: $O(\log n)$ Modify: $O(\log n)$
- Hash table
 - Insert: $O(1)$ Lookup: $O(1)$ Modify: $O(1)$

Hash Tables

- Custom data structure for fast key/value lookups
 - Used to implement the Map ADT
 - Goal: $O(1)$ access (insert/modify/delete)
 - So we can't use comparisons (remember sorting?)
- Observation: arrays provide $O(1)$ access
 - How to map from keys to array indices?
 - How large does the array need to be?



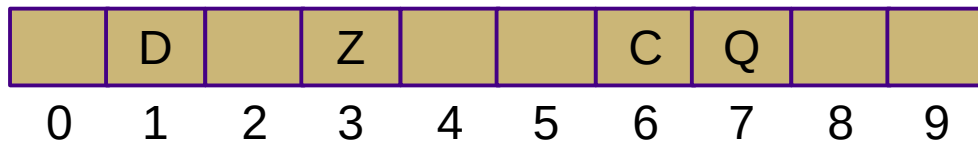
Map ADT



Array

Hash Tables

- Simple case: keys are integers in $[0, N)$
 - Create an array of length N
 - Use keys directly as indices into the array
- This does not scale!
 - N could be very large
 - Keys might not be integers

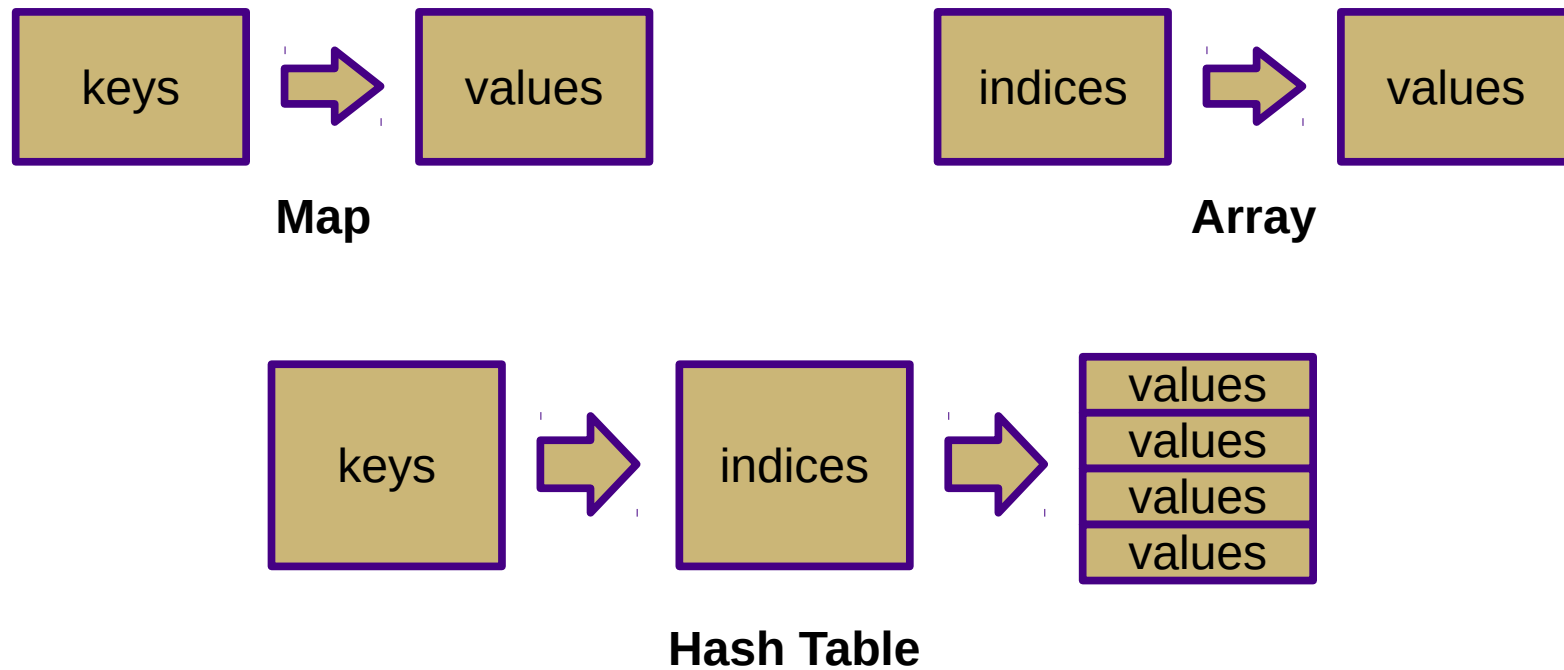


Items: (1,D), (3,Z), (6,C), (7,Q)

```
char t[10];  
t[3] = 'Z';  
t[6] = 'C';  
t[7] = 'Q';  
t[1] = 'D';
```

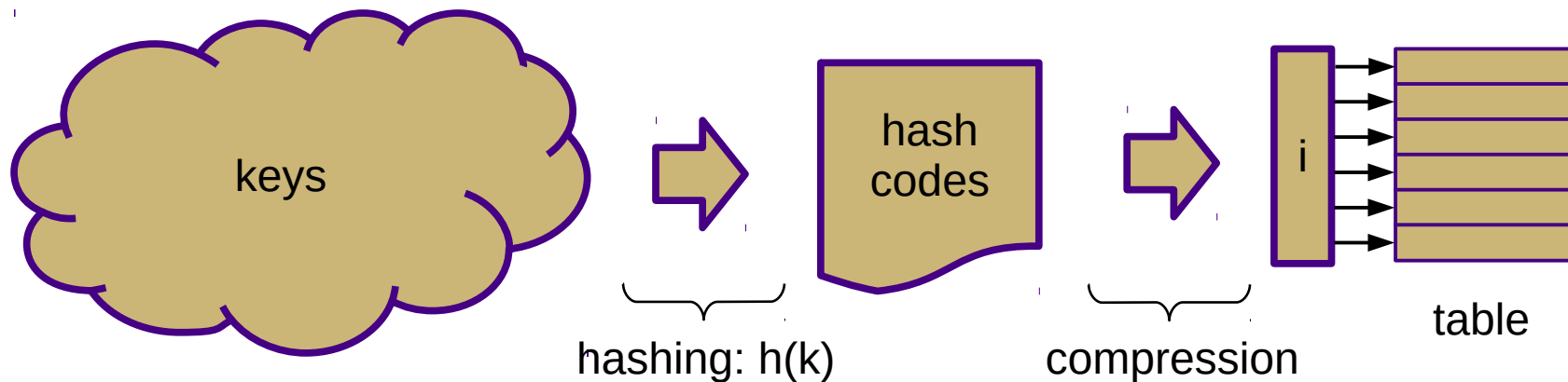
Hash Tables

- Main concept: hash function
 - Converts keys → table indices
 - Table holds "buckets" of elements



Hash Functions

- Hash code (key \rightarrow 32/64-bit integer)
 - Translation from key domain to hash code domain
 - Key can be any immutable object
 - Hash codes are usually native integers
- Compression function (hash code \rightarrow table index)
 - Compression from hash code domain to index domain
 - Result is used to access table storage



Hash Codes

- Most codes are based on interpreting raw bits as integers
 - Issue: key size may be greater than native integer width
 - Need to combine multiple integers
 - Truncation
 - Summation
 - Exclusive-or (XOR)
 - Polynomial combination
 - Cyclic shifting
 - Cryptographic hashes (e.g., MD5, SHA-1)
- } bad for variable-length objects

Compression Functions

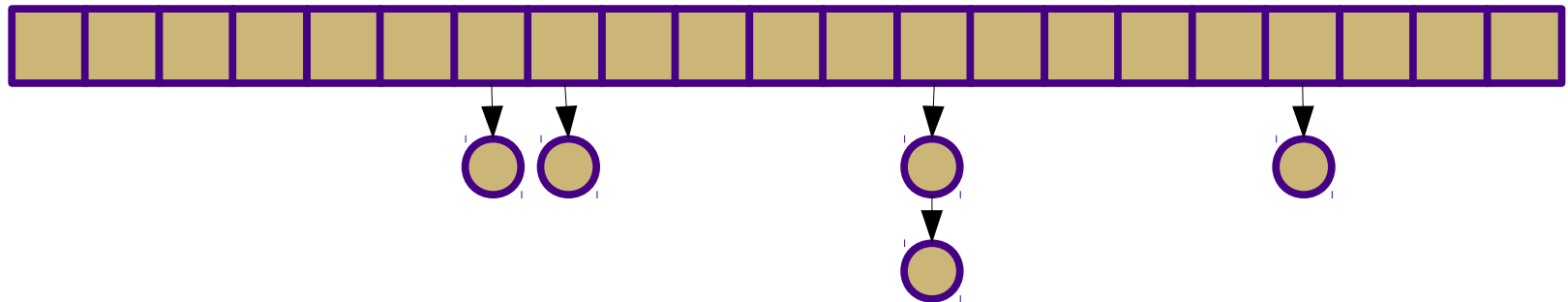
- Simplest: Modulus division
 - $h' = h(k) \% N$
 - N is the number of buckets in the hash table
 - N should be a prime number
- Better: Multiply-Add-and-Divide
 - $h' = ((a \cdot h(k)) + b) \% p) \% N$
 - p is a prime number larger than N
 - a and b are random integers from $[0, p-1]$
 - $a > 0$
 - Essentially a pseudo-random number generator that uses hash codes as seeds

Hash Functions

- Major problem: **collisions**
 - Multiple keys that map to the same index
 - Two-fold approach:
 - Minimize collisions by choosing a good hash code
 - Handle collisions with chaining or probing

Collision Handling

- Separate chaining
 - Each bucket is a linked list of elements
 - Load factor: $\lambda = n/N$
 - Expected size of each bucket
 - If the hash function is good, map operations run in $O(\lambda)$
 - This should be a small constant
 - Preferably less than 1
 - As long as λ is $O(1)$, map operations run in $O(1)$ expected time

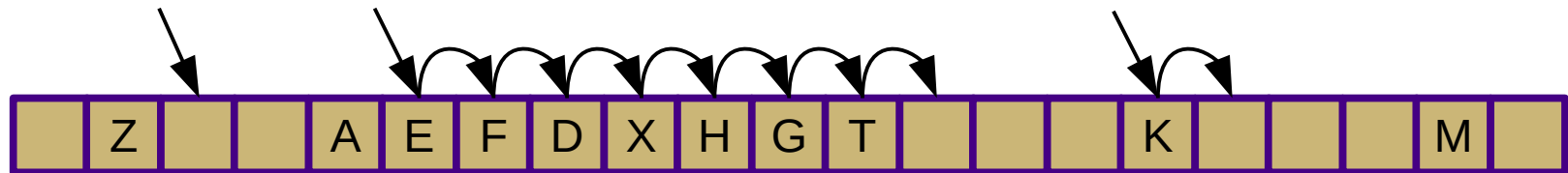


Collision Handling

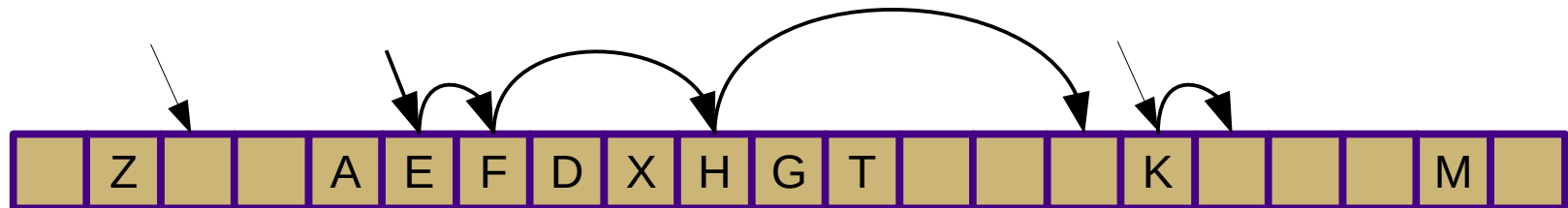
- Open addressing / probing
 - Only one (key, value) pair per "bucket"
 - Problem: $h'(k)$ not guaranteed to be open
 - Probing scheme to find an open bucket
 - Load factor: $\lambda = n/N$
 - Percentage of buckets that are occupied
 - Allowing removal of keys from a map can break this scheme!
- Approaches
 - Linear probing: $(h'(k) + i) \% N$
 - Quadratic probing: $(h'(k) + i^2) \% N$
 - Double hashing: $(h'(k) + i \cdot h''(k)) \% N$
 - Pseudo-random probing: $(h'(k) + \text{prand}(i)) \% N$

Open Addressing

- Linear probing

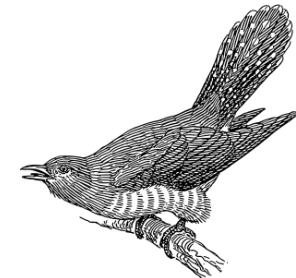


- Quadratic probing



Collision Handling

- Coalesced hashing (hybrid chained/open)
 - Maintain chains as pointers between buckets
 - Avoids some of the overhead of probing
- Cuckoo hashing (multiple hash functions)
 - Use multiple hash functions (primary and alternate)
 - If new key's bucket is full, remove existing key and re-insert it using alternate hash
 - Repeat until empty bucket is found or an infinite loop is detected



Load Factors

- Separate chaining
 - Want to keep λ less than 1 (preferably < 0.9)
- Open addressing
 - Want to keep λ less than $1/2$ or $2/3$
- Rehashing
 - When constraints above are violated, resize the hash table and re-apply the compression function to re-insert all keys
 - Cost can be amortized by doubling the table size
 - Just as with dynamic arrays

Hashing Analysis

- The expected # of keys in a bucket is $\text{ceil}(n/N)$
 - This is $O(1)$ if n is $O(N)$
 - Assumes a good hash function
 - Assumes enforcement of appropriate load factor
- Thus, expected costs for major map operations (insertion, modification, lookup, removal) are all $O(1)$
 - Worse case: $O(n)$
- Full probabilistic analysis is beyond the scope of this class

Retrospective

- Set/Map equivalence
 - A Set is a Map with no values
 - A Map is a Set of keys with associated values
 - In general, any implementation of one can be used for the other
- Progression of Set/Map implementations:
 - Array / linked list
 - Mostly $O(n)$ operations
 - Skip list / AVL tree
 - Mostly $O(\log n)$ operations
 - Hash table
 - Mostly $O(1)$ operations