Hash Tables continued Prefix Trees

CS 311 Data Structures and Algorithms Lecture Slides Wednesday, November 18, 2020

Glenn G. Chappell Department of Computer Science University of Alaska Fairbanks ggchappell@alaska.edu © 2005-2020 Glenn G. Chappell Some material contributed by Chris Hartman

Review

- Our problem for most of the rest of the semester:
 - Store: A collection of data items, all of the same type.
 - Operations:
 - Access items [single item: retrieve/find, all items: traverse].
 - Add new item [insert].
 - Eliminate existing item [delete].
 - Time & space efficiency are desirable.

A solution to this problem is a **container**.

In a generic container, client code can specify the value type.

Unit Overview Tables & Priority Queues

Majo	r Topics		
✓ ∎	Introduction to Tables <	- Lots of lousy implementations	
✓ ■	Priority Queues		
✓ ■	Binary Heap Algorithms	Idea #1: Restricted Table	
✓ ■	Heaps & Priority Queues in the C++ STL		
✓ ∎	2-3 Trees	Idaa #2. Kaap a traa balancad	
✓ ■	Other self-balancing search trees		
(part) 🛛	Hash Tables	Idea #3: Magic functions	
	Prefix Trees	- A special-purpose	
1.1	Tables in the C++ STL & Elsewhere	implementation: "the Radix Sort of Table implementations"	

A **Table** allows for arbitrary key-based look-up.
Three single-item operations: retrieve, insert, delete by key.
A Table implementation typically holds **key-value pairs**.



Three ideas for efficient implementations:

- 1. Restricted Table \rightarrow Priority Queues
- 2. Keep a tree balanced \rightarrow Self-balancing search trees
- 3. Magic functions \rightarrow Hash Tables

Overview of Advanced Table Implementations

We cover the following advanced Table implementations.

- Self-balancing search trees
 - To make things easier, allow more children (?):
 - ✓ 2-3 Tree
 - Up to 3 children
 - ✓ 2-3-4 Tree
 - Up to 4 children
 - Red-Black Tree
 - Binary Tree representation of a 2-3-4 Tree
 - Or back up and try for a strongly balanced Binary Search Tree again:
 - AVL Tree
- Alternatively, forget about trees entirely:

(part) Hash Table

- Finally, "the Radix Sort of Table implementations":
 - Prefix Tree

Idea #2: Keep a tree balanced

Later, we will cover other self-balancing search trees: B-Trees, B+ Trees.

Idea #3: Magic functions The **self-balancing search trees** include 2-3 Trees, 2-3-4 Trees, Red-Black Trees, AVL Trees, and other kinds.

- All are generalizations of Binary Search Trees.
- The ones we covered all allow for Table implementations with logarithmic-time retrieve/insert/delete by key.



Generally, the Red-Black Tree is agreed to have the best *overall* performance, for in-memory datasets with many insert & delete operations, when worst-case performance is important.

2020-11-18

A **Hash Table** is a Table implementation that stores <u>key-value</u> pairs in an unsorted array. Array indices are **slots**.

- A key's slot is computed using a **hash function**.
- An array location can be *EMPTY*.

Or just keys, if there are no associated values.



Collision: when two items with unequal keys get the same slot. Collisions are generally an unavoidable problem, because there are often far more possible keys than slots.

Needed

- Hash function (typically separate from the Hash Table implementation).
- Collision-resolution method.



A hash function *must*:

- Take a key and return a nonnegative integer (hash code).
- Be deterministic: the output depends only on the input. Passing the same key multiple times results in the same hash code every time.
- Return the same hash code for keys that compare equal (==).

Consistency requirement mentioned previously

A *good* hash function:

- Is fast.
- Spreads out its results evenly over the possible output values.
- Turns patterns in its input into unpatterned output.

2020-11-18

Collision resolution methods, category #1: **Closed Hashing**

- Each array item holds one key-value pair, or a mark indicating EMPTY or DELETED.
- To find a key, begin at the slot given by the hash function, and probe in a sequence of slots: the probe sequence. End when desired key or EMPTY slot is found.

Some probe sequences (*t* is initial slot given by hash function):

- **Linear probing**: *t*, *t*+1, *t*+2, *t*+3, etc.
 - Tends to form **clusters**, which slow things down. \otimes
- **Quadratic probing**: *t*, *t*+1², *t*+2², *t*+3², etc.
- Double hashing: Base probe sequence on a 2nd hash function.



10

Collision resolution methods, category #2: **Open Hashing**

- Each array item holds a data structure (a **bucket**) that can store multiple key-value pairs.
- Buckets are virtually always Singly Linked Lists.
- To find a key, determine which bucket to look in based on the hash code. Do a Sequential Search in that bucket.



Hash Tables

continued

Hash Tables Rehashing

- All Hash Table implementations that allow insertion will suffer poor performance when the data structure gets too full.
- Q. What do we do about this?
- A. When the number of items gets too high, we remake the Hash Table, doing a reallocate-and-copy to a larger array—as we did with resizable smart arrays. This is called **rehashing**.

Rehashing is time-consuming. We need to traverse the entire Hash Table, calling the hash function for every key present. This is one of the downsides of Hash Tables.

Two Questions

1. How do we decide when to do rehashing?

2. How does periodic rehashing affect Hash Table performance? We look at these after discussing basic Hash Table efficiency. How efficient are Hash Tables?

Let's look at Hash Tables where duplicate keys are not allowed:

- A Hash Table using open hashing, in which buckets are Linked Lists.
- A Hash Table using closed hashing.



For now, assume there is no rehashing—a dumb assumption, but only a temporary one.

The time taken for *retrieve*, *insert*—assuming no rehashing—or *delete* is essentially the time required to **search** by key (right?).

So: how fast can we do a search by key?

2020-11-18

In the worst case, every item inserted gets the same slot.

With open hashing, this turns our Table into a Linked List. So a search is linear-time.

With closed hashing, a search may require probing every stored item. Again, linear-time.



Conclusion. In a Hash Table, search by key is linear-time. Thus, *retrieve*, *insert*, and *delete* are linear-time (worst case!).

Now, what about the average case?

2020-11-18

Hash Tables Efficiency — Load Factor [1/2]

Average-case performance of a Hash Table can be analyzed based on its **load factor**: $\alpha = (\# \text{ of keys present}) / (\# \text{ of slots})$.

> Lower-case Greek letter **alpha**





Answers on next slide.

Average-case performance of a Hash Table can be analyzed based on its **load factor**: $\alpha = (\# \text{ of keys present}) / (\# \text{ of slots})$.



A Hash Table implementation keeps its load factor small.

In the following slides, we assume α is significantly less than 1. For example, we might require that $\alpha < 2/3$.

Consider open hashing. Again, buckets = Linked Lists, no duplicate keys, α significantly less than 1.

Search worst case: linear time.



Search average case:

- The average number of items in a bucket is α (the load factor).
- Thus, the average number of comparisons required for a search resulting in NOT FOUND is α.
- The average number of comparisons required for a search resulting in FOUND is approximately 1 + α/2.
 - Average search looks at item found + half of other items in bucket:

$$1 + \frac{1}{2} \cdot \frac{n-1}{b} = 1 + \frac{n/b}{2} - \frac{1}{2b} \approx 1 + \frac{\alpha}{2}$$
Very
small

If $\alpha < 1$, then these are less than 1.5. That's *fast*!

```
2020-11-18
```

Consider closed hashing. Again, no duplicate keys, α significantly less than 1.



Search average case:

- Comparisons for linear probing:
 - NOT FOUND: (1/2)[1+1/(1-α)]².
 - FOUND: (1/2)[1+1/(1-α)].
- Comparisons for quadratic probing:
 - NOT FOUND: 1/(1-α).
 - FOUND: -ln(1-α)/ α.

Important to know:

- Larger α means a slower average case for search.
- Requiring α < r, for r a fixed number between 0 and 1 (so r ≠ 1), gives a fast constant-time average case for search.

The analysis is complicated. It actually took several years to be properly figured out. We will not give details.

2020-11-18

Summary, so far. For both open & closed hashing:

- Worst-case performance for *retrieve*, *insert*, *delete*: linear time.
- Average-case performance for *retrieve*, *delete*: constant time. Also for *insert*—not counting the time required for rehashing.

But of course we will need to do rehashing occasionally.

1. How do we decide when to do rehashing?

The Two Questions from a few slides back

- Primary trigger for rehashing: the load factor becomes too large.
- 2. How does periodic rehashing affect Hash Table performance?
 - The effect of rehashing on Hash-Table efficiency is similar to that of reallocate-and-copy on a resizable array. Again, when we increase the array size, we need to increase it by a *constant factor*.
 - When we count rehashing, *insert* in a well written Hash Table will have an average case of amortized constant time.
 - Note the *double average*! Average data & amortized time.

Efficiency Comparison (duplicate keys not allowed)

	Idea #1	Idea #2	Idea #3	
	Priority Queue using Heap	Self-Balancing Search Tree	Hash Table: worst case	Hash Table: average case
Retrieve	Constant*	Logarithmic	Linear	Constant
Insert	Amortized** logarithmic	Logarithmic	Linear	Amortized constant***
Delete	Logarithmic*	Logarithmic	Linear	Constant

*Priority Queue *retrieve* & *delete* are not Table operations in full generality. Only the item with the highest priority (key) can be retrieved/deleted.

- **Logarithmic if enough memory is preallocated. Otherwise, occasional reallocateand-copy—linear time—may be required. Time per *insert*, averaged over many consecutive *inserts*, will be logarithmic. Thus, *amortized logarithmic time* (which is not a term I expect you to know).
- ***Hash Table *insert* is constant-time only in a *double average* sense: averaged both over all possible inputs and over a large number of consecutive *inserts*.

Hash Tables generally have excellent average-case performance, but poor worst-case performance.

Can a **malicious user** force poor Hash Table performance?

- Sometimes. For many applications this matters little, but for others it can matter a great deal. This issue is something to keep in mind.
- A cryptographic hash function is one that is extremely difficult to reverse. This makes it hard to force poor performance. However, cryptographic hash functions are time-consuming to compute, which makes them unsuitable for most applications of Hash Tables.

Hash Tables are appropriate for many use cases of Tables. But their worst-case performance can be problematic in some contexts.

When using Hash Tables, do so intelligently.

Don't program

a pacemaker

in Python!

Prefix Trees

Prefix Trees Background

In a list of strings, there are often strings that start with the same sequence of characters.

For example, in the list to the right, "dot", "dote", and "doting" all begin with "dot...".

Such strings are said to have a common **prefix**.

We will use *string* in a general sense, just as when we discussed Radix Sort. A **string** is a sequence; its entries will be called **characters**.

- The words in the above list are strings of letters.
- A nonnegative integer can be regarded as a string of digits.
- Many kinds of data can be regarded as strings of **bits** (0s & 1s).

We look at a data structure that can be used to store a Table in which the keys are such strings: a *Prefix Tree*.



Prefix Trees Definition

- A **Prefix Tree** (a.k.a. **Trie**) is a tree used to implement a Table whose keys are strings—in the general sense.
 - Each child is associated with a character.
 A node has at most one child for each character.
 - A node has:
 - A Boolean—whether it represents a stored key.
 - Child pointers—one for each possible character.
 - The value associated with a key, if needed.
 - Duplicate keys are generally not allowed. (But they can be faked; do you see how?)

Nodes with dots represent stored keys: **dig**, **dog**, **dot**, **dote**, **doting**, **eggs**.

Credit

- René de la Briandais 1959, Edward Fredkin 1960.
- "Trie" was coined by Fredkin.

"**Trie**" for "re**TRIE**val". Some say "tree". Some say "try". I say "Prefix Tree".



Prefix Trees Efficiency

For a Prefix Tree, *retrieve*, *insert*, and *delete* by key all take a number of steps proportional to the **length of a key**.If key length is considered fixed, then all are constant time!

But larger datasets tend to have longer keys.

If C is the size of the character set, and L is the length of a key, then the number of possible keys is $K = C^{L}$.

Solve for L, and we get $L = \log_C K$.

The **length of a key** is logarithmic in the number of possible keys.

So there is a hidden logarithm, just like Radix Sort.



Prefix Trees Thoughts

A Prefix Tree is a good basis for a Table implementation, when keys are short-ish sequences of characters from a small-ish set.

- Words in a dictionary, ZIP codes, etc.
- Just like Radix Sort.

Prefix Trees are **easy to implement well**.

- If you feel like writing a Red-Black Tree to be used in production code, then you might want to sit down until the feeling goes away.
- But if you feel like writing a Prefix Tree, then go for it!

The idea behind Prefix Trees is also used in other data structures. We will look at one of these later.

