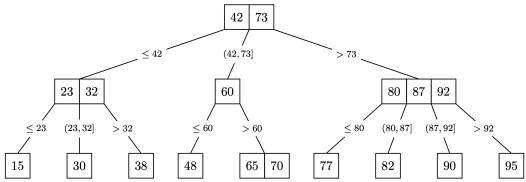
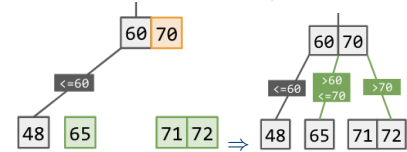


(a, b)-trees

e.g. a (2, 4)-tree storing 18 keys



- rules
 - (a, b) -child policy where $2 \leq a \leq (b + 1)/2$
- an internal node has 1 more child than its number of keys
- all leaf nodes must be at the same depth from the root
- terminology (for a node z)
 - key range - range of keys covered in subtree rooted at z
 - keylist - list of keys within z
 - treelist - list of z 's children
- max height = $O(\log_a n) + 1$
- min height = $O(\log_b n)$
- $\text{search}(\text{key}) \Rightarrow O(\log n)$
 - $= O(\log_2 b \cdot \log_a n)$ for binary search at each node
- $\text{insert}(\text{key}) \Rightarrow O(\log n)$
- $\text{split}()$ a node with too many children
 - use median to split the keylist into 2 halves
 - move median key to parent; re-connect remaining nodes
 - (if the parent is now unbalanced, recurse upwards; if the root is reached, median key becomes the new root)



- $\text{delete}(\text{key}) \Rightarrow O(\log n)$
 - if the node becomes empty, $\text{merge}(y, z)$ - join it with its left sibling & replace it with their parent
- if the combined nodes exceed max size: $\text{share}(y, z) = \text{merge}(y, z)$ then $\text{split}()$

B-Tree

- $(B, 2B)$ -trees $\Rightarrow (a, b)$ -tree where $a = B, b = 2B$
- possible augmentation: use a LinkedList to connect between each level

Merkle Trees

- binary tree - nodes augmented with a hash of their children
- same root value = identical tree

HASH TABLES

- disadvantage: no successor/predecessor operation

hashing

- Let the m be the table size; let n be the number of items; let $\text{cost}(h)$ be the cost of the hash function
- $\text{load}(\text{hash table}), \alpha = \frac{n}{m}$
 - = average number of items per bucket
 - = expected number of items per bucket

hashing assumptions

- simple uniform hashing assumption
 - every key has an equal probability of being mapped to every bucket
 - keys are mapped independently
- uniform hashing assumption
 - every key is equally likely to be mapped to every permutation, independent of every other key.
 - NOT fulfilled by linear probing

properties of a good hash function

- able to enumerate all possible buckets - $h : U \rightarrow \{1..m\}$
 - for every bucket $j, \exists i$ such that $h(\text{key}, i) = j$
- simple uniform hashing assumption

hashCode

rules for the hashCode() method

- always returns the same value, if the object hasn't changed
- if two objects are equal, they return the same hashCode

rules for the equals method

- reflexive - $x.\text{equals}(x) \Rightarrow \text{true}$
- symmetric - $x.\text{equals}(y) \Rightarrow y.\text{equals}(x)$
- transitive - $x.\text{equals}(y), y.\text{equals}(z) \Rightarrow x.\text{equals}(z)$
- consistent - always returns the same answer
- null is null - $x.\text{equals}(\text{null}) \Rightarrow \text{false}$

chaining

- time complexity
 - $\text{insert}(\text{key}, \text{value}) - O(1 + \text{cost}(h)) \Rightarrow O(1)$
 - for n items: expected maximum cost
 - $= O(\log n)$
 - $= \Theta(\frac{\log n}{\log(\log(n))})$
 - $\text{search}(\text{key})$
 - worst case: $O(n + \text{cost}(h)) \Rightarrow O(n)$
 - expected case: $O(\frac{n}{m} + \text{cost}(h)) \Rightarrow O(1)$
- total space: $O(m + n)$

open addressing - linear probing

- redefined hash function: $h(k, i) = h(k, 1) + i \bmod m$
- $\text{delete}(\text{key})$
 - use a tombstone value - DON'T set to null
- performance
 - if the table is $\frac{1}{4}$ full, there will be clusters of size $\Theta(\log n)$
 - expected cost of an operation, $E[\#probes] \leq \frac{1}{1-\alpha}$ (assume $\alpha < 1$ and uniform hashing)
- advantages
 - saves space (use empty slots vs linked list)

- better cache performance (table is one place in memory)
- rarely allocate memory (no new list-node allocation)
- disadvantages
 - more sensitive to choice of hash function (clustering)
 - more sensitive to load (as $\alpha \rightarrow 1$, performance degrades)

double hashing

- for 2 functions f, g , define $h(k, i) = f(k) + i \cdot g(k) \bmod m$
- if $g(k)$ is relatively prime to m , then $h(k, i)$ hits all buckets
 - e.g. for $g(k) = n^k, n$ and m should be coprime.

table size

- assume chaining & simple uniform hashing
- let m_1 = size of the old hash table; m_2 = size of the new hash table; n = number of elements in the hash table
- growing the table: $O(m_1 + m_2 + n)$
 - rate of growth

table growth	resize	insert n items
increment by 1	$O(n)$	$O(n^2)$
double	$O(n)$	$O(n)$, average $O(1)$
square	$O(n^2)$	$O(n)$

PROBABILITY THEORY

- if an event occurs with probability p , the expected number of iterations needed for this event to occur is $\frac{1}{p}$.
- for random variables: expectation is always equal to the probability
- linearity of expectation: $E[A + B] = E[A] + E[B]$

UNIFORMLY RANDOM PERMUTATION

- for an array of n items, every of the $n!$ possible permutations are producible with probability of exactly $\frac{1}{n!}$
 - the number of outcomes should distribute over each permutation uniformly. (i.e. $\frac{\# \text{ of outcomes}}{\# \text{ of permutations}} \in \mathbb{N}$)
- probability of an item remaining in its initial position = $\frac{1}{n}$
- KnuthShuffle $\Rightarrow O(n)$ - for every element in array A , swap it with a random index in array A .

sort	best	average	worst	stable?	memory
bubble	$\Omega(n)$	$O(n^2)$	$O(n^2)$	✓	$O(1)$
selection	$\Omega(n^2)$	$O(n^2)$	$O(n^2)$	×	$O(1)$
insertion	$\Omega(n)$	$O(n^2)$	$O(n^2)$	✓	$O(1)$
merge	$\Omega(n \log n)$	$O(n \log n)$	$O(n \log n)$	✓	$O(n)$
quick	$\Omega(n \log n)$	$O(n \log n)$	$O(n^2)$	×	?

sorting invariants		searching	
sort	invariant (after k iterations)	search	average
bubble	largest k elements are sorted	linear	$O(n)$
selection	smallest k elements are sorted	binary	$O(\log n)$
insertion	first k elements are in order	quickSelect	$O(n)$
merge	—	interval	$O(\log n)$
quick	partition is in the right position	all-overlaps	$O(k \log n)$
		1D range	$O(k + \log n)$
		2D range	$O(k + \log^2 n)$

data structures assuming $O(1)$ comparison cost		
data structure	search	insert
sorted array	$O(\log n)$	$O(n)$
unsorted array	$O(n)$	$O(1)$
linked list	$O(n)$	$O(1)$
tree (kd/(a, b)/binary)	$O(\log n)$ or $O(h)$	$O(\log n)$ or $O(h)$
trie	$O(L)$	$O(L)$
dictionary	$O(\log n)$	$O(\log n)$
symbol table	$O(1)$	$O(1)$
chaining	$O(n)$	$O(1)$
open addressing	$\frac{1}{1-\alpha} = O(1)$	$O(1)$

orders of growth	
$T(n) = 2T(\frac{n}{2}) + O(n)$	$\Rightarrow O(n \log n)$
$T(n) = T(\frac{n}{2}) + O(n)$	$\Rightarrow O(n)$
$T(n) = 2T(\frac{n}{2}) + O(1)$	$\Rightarrow O(n)$
$T(n) = T(\frac{n}{2}) + O(1)$	$\Rightarrow O(\log n)$
$T(n) = 2T(n - 1) + O(1)$	$\Rightarrow O(2^n)$
$T(n) = 2T(\frac{n}{2}) + O(n \log n)$	$\Rightarrow O(n(\log n)^2)$
$T(n) = 2T(\frac{n}{4}) + O(1)$	$\Rightarrow O(\sqrt{n})$
$T(n) = T(n - c) + O(n)$	$\Rightarrow O(n^2)$